

1 **Cohort profile: St. Michael's Hospital Tuberculosis Database (SMH-TB), a retrospective**
2 **cohort of electronic health record data and variables extracted using natural language**
3 **processing**

4 David Landsman¹, Ahmed Abdelbasit², Christine Wang², Michael Guerzhoy^{3, 4, 5}, Ujash Joshi⁴,
5 Shaun Mathew⁶, Chloe Pou-Prom⁷, David Dai⁷, Victoria Pequegnat⁸, Joshua Murray⁷, Kamalprit
6 Chokar⁹, Michaelia Banning⁷, Muhammad Mamdani^{5, 7, 10, 11, 12}, Sharmistha Mishra^{¶, 1, 2, 11}, Jane
7 Batt^{*,¶, 13, 14}

- 8 1. MAP Centre for Urban Health Solutions, Li Ka Shing Knowledge Institute, St. Michael's
9 Hospital, Unity Health Toronto, Toronto, Ontario, Canada
10 2. Department of Medicine, University of Toronto, Toronto, Ontario, Canada
11 3. Princeton University, Princeton, New Jersey, United States
12 4. University of Toronto, Toronto, Ontario, Canada
13 5. Li Ka Shing Knowledge Institute, St. Michael's Hospital, Unity Health Toronto, Toronto,
14 Ontario, Canada
15 6. Department of Computer Science, Ryerson University, Toronto, Ontario, Canada
16 7. Unity Health Toronto, Toronto, Ontario, Canada
17 8. Decision Support Services, St. Michael's Hospital, Unity Health Toronto, Toronto,
18 Ontario, Canada
19 9. Division of Respiriology, Department of Medicine, St. Michael's Hospital, Unity Health
20 Toronto, Toronto, Ontario, Canada
21 10. Leslie Dan Faculty of Pharmacy, University of Toronto, Canada, Toronto, Ontario,
22 Canada
23 11. Institute of Health Policy, Management, and Evaluation, University of Toronto, Toronto,
24 Ontario, Canada
25 12. Vector Institute, Toronto, Ontario, Canada
26 13. Keenan Research Center for Biomedical Science, St. Michael's Hospital, Unity Health
27 Toronto, Toronto, Ontario, Canada
28 14. Institute of Medical Science, University of Toronto, Toronto, Ontario, Canada

29 [¶]These authors contributed equally to this work.

30 ^{*}Corresponding author: Jane Batt, MD, PhD

31 E-mail: Jane.batt@utoronto.ca

32

33

34

35 **Author Contributions**

36 David Landsman – Data Curation, Formal Analysis, Investigation, Methodology, Software,
37 Validation, Visualization, Writing – Original Draft Preparation, Writing – Review & Editing

38 Ahmed Abdelbasit – Data Curation, Investigation, Validation, Writing – Original Draft
39 Preparation, Writing – Review & Editing

40 Christine Wang – Data Curation, Investigation, Validation, Writing – Review & Editing

41 Michael Guerzhoy – Investigation, Methodology, Software, Writing – Review & Editing

42 Ujash Joshi – Investigation, Methodology, Software, Writing – Review & Editing

43 Shaun Mathew – Investigation, Methodology, Software, Writing – Review & Editing

44 Chloe Pou-Prom – Investigation, Methodology, Software, Writing – Review & Editing

45 David Dai – Investigation, Methodology, Software, Writing – Review & Editing

46 Victoria Pequegnat – Data Curation, Resources, Validation, Writing – Review & Editing

47 Joshua Murray – Investigation, Methodology, Software, Writing – Review & Editing

48 Kamalprit Chokar – Data Curation, Writing – Review & Editing

49 Michaelia Banning – Project Administration, Writing – Review & Editing

50 Muhammad Mamdani – Project Administration, Writing – Review & Editing

51 Sharmistha Mishra – Conceptualization, Data Curation, Funding Acquisition, Investigation,
52 Methodology, Project Administration, Resources, Supervision, Validation, Writing – Review &
53 Editing

54 Jane Batt – Conceptualization, Data Curation, Funding Acquisition, Investigation, Methodology,
55 Project Administration, Resources, Supervision, Validation, Writing – Review & Editing

56 **Abstract**

57 ***Background***

58 Tuberculosis (TB) is a major cause of death worldwide. TB research draws heavily on clinical
59 cohorts which can be generated using electronic health records (EHR), but granular information
60 extracted from unstructured EHR data is limited. The St. Michael's Hospital TB database (SMH-
61 TB) was established to address gaps in EHR-derived TB clinical cohorts and provide researchers
62 and clinicians with detailed, granular data related to TB management and treatment.

63 ***Methods***

64 We collected and validated multiple layers of EHR data from the TB outpatient clinic at St.
65 Michael's Hospital, Toronto, Ontario, Canada to generate the SMH-TB database. SMH-TB
66 contains structured data directly from the EHR, and variables generated using natural language
67 processing (NLP) by extracting relevant information from free-text within clinic, radiology, and
68 other notes. NLP performance was assessed using recall, precision and F_1 score averaged across
69 variable labels. We present characteristics of the cohort population using binomial proportions
70 and 95% confidence intervals (CI), with and without adjusting for NLP misclassification errors.

71 ***Results***

72 SMH-TB currently contains retrospective patient data spanning 2011 to 2018, for a total of 3298
73 patients (N=3237 with at least 1 associated dictation). Performance of TB diagnosis and
74 medication NLP rulesets surpasses 93% in recall, precision and F_1 metrics, indicating good
75 generalizability. We estimated 20% (95% CI: 18.4-21.2%) were diagnosed with active TB and
76 46% (95% CI: 43.8-47.2%) were diagnosed with latent TB. After adjusting for potential
77 misclassification, the proportion of patients diagnosed with active and latent TB was 18% (95%
78 CI: 16.8-19.7%) and 40% (95% CI: 37.8-41.6%) respectively

79 ***Conclusion***

80 SMH-TB is a unique database that includes a breadth of structured data derived from structured
81 and unstructured EHR data. The data are available for a variety of research applications, such as
82 clinical epidemiology, quality improvement and mathematical modelling studies.

83

84 **Introduction**

85 Tuberculosis (TB) is the top infectious killer worldwide, resulting in 1.6 million deaths in 2017
86 (1). 1.7 billion people carry the latent form of the infection, of whom 10% at minimum, will
87 develop the active, infectious form of disease. Latent TB infection (LTBI) progression to active
88 disease can be prevented and TB can be cured, with appropriate antibiotics taken over many
89 months. TB is endemic in many low-income countries and particularly prevalent in Asia and
90 Africa. The World Health Organization recommends the treatment of LTBI as part of the global
91 “End TB Strategy”, and an achievable goal critical to TB elimination in high-income countries
92 (2,3).

93 Given the burden of active TB disease is disproportionately carried in low-resource settings,
94 research addressing disease epidemiology, treatment (including clinical trials and programs of
95 delivery), and the use and utility of innovative and point of care diagnostics is often completed in
96 the populations of countries with highest burden of TB. The prevalence of LTBI on the other
97 hand, is considerable even in high-income countries (CDC estimates 13,000,000 people living
98 the USA have LTBI (4)) and thus research ranging from basic pathogenesis to program
99 development can be conducted on the global population. Indeed while advances in biomedical
100 research over the past 1 to 2 decades have delivered successes ranging from rapid point-of-care
101 diagnostics testing for pulmonary TB to the development of novel therapeutics such as
102 bedaquiline and delamanid, many questions remain, including, for example, discovering
103 biomarkers that precisely indicate individuals at risk of LTBI activation and developing
104 programs of TB care that ensure efficacy, are equitable and resilient (1,5).

105 Many primary care practices and hospitals in high-income countries have curated electronic
106 health record (EHR) data for research and surveillance (6–9), that improve ease of access to
107 information and data sharing for collaborative work. The use of EHRs in hospital and office-
108 based clinical practices has risen substantially in the past decade, providing rich data sources that
109 have the potential to simultaneously improve patient care and advance research initiatives
110 (10,11). Most EHR-derived databases are however limited to structured data, such as
111 demographic information collected at patient registration, laboratory tests and results and
112 diagnostic codes used in physician billing. As such, the rich, granular data embedded within
113 unstructured (text) data from dictated notes on both hospital admitted and clinic patients are
114 excluded (12,13) unless these variables are abstracted via manual chart review (14,15) or natural
115 language processing (NLP) (16–18). We developed the first digital retrospective clinical
116 database that combines structured data, unstructured (text) data, and variables derived from
117 transforming unstructured data to structured data using natural language rulesets, among patients
118 assessed in an inner-city outpatient TB clinic at St Michaels Hospital (SMH) of Unity Health
119 Toronto in Toronto, Ontario, Canada.

120 Approximately 2000 people (5.6 per 100,000 people) are diagnosed with active TB in Canada
121 (19) annually and 1.3 million are estimated to have LTBI. The SMH TB clinic cares exclusively
122 for individuals with suspected or diagnosed active TB and LTBI, seeing 1800-2200 patient-visits
123 each year, and assessing and developing a diagnostic and management plan for 670-800 new
124 patients each year.

125 In this paper we describe the SMH-TB database, which aims to be a resource for scientists who
126 are conducting research into many facets of TB, ranging from observational epidemiology to
127 emulated trials and quality improvement and implementation science research. The purpose of
128 this profile is to describe our methodology, present the cohort and the database validation.
129 Access to the database is available to collaborators wishing to work with the research team of the
130 SMH TB clinic. The NLP rulesets developed to extract variables from the unstructured data in
131 the EHR are publicly available on GitHub (20).

132 **Materials and methods**

133 **Cohort Description**

134 The database compiles all data available on all TB clinic patients (N=3298) treated at SMH from
135 April 2011 to December 2018. The database contains socio-demographic information
136 surrounding immigration, housing status, and insurance, and clinical information including
137 laboratory and imaging results, co-morbidities, diagnoses and treatment. Ethics approval for
138 development and validation of the database was obtained from the Unity Health Toronto
139 Research Ethics Board (REB 19-080). Patient consent was not required or obtained as per the
140 Tri-Council Policy Statement 2 (TCPS2), since only retrospective data were collected from
141 clinical charts (21).

142 Patients are referred to the TB outpatient clinic predominantly from Public Health Units in the
143 Greater Toronto area (population of 6 million), Canada Immigration and Citizenship,
144 Occupational Health and Safety Departments of Toronto area hospitals, community health care
145 professionals (physicians, nurse-practitioners), and SMH staff physicians caring for an admitted
146 patient or a patient in the emergency room (ER). When including a patient in our database we
147 consider all available encounters, including inpatient admissions and ER visits.

148 **Data Collection**

149 St. Michael's Hospital EHR is managed by several systems. The Enterprise Data Warehouse
150 (EDW) stores and manages structured data including patient demographics and medical test
151 results. Soarian stores the unstructured patient data, which includes dictated clinical notes. SMH-
152 TB retrieved data of patients registered and assessed in the TB outpatient clinic to provide a
153 comprehensive description of patient characteristics, disease, management and clinical trajectory.

154 SMH-TB is restricted to a start-date of April 2011, which is the date of initiation of EHR at
 155 SMH. Fig 1 shows the data flow and data sources for the SMH-TB database.

156 **Fig 1: Data sources for SMH-TB Database**

157 The SMH-TB database stores patient characteristics and encounter data in separate tables, which
 158 can be linked together using unique, de-identified patient or encounter IDs. Fig 2 presents the
 159 tables provided in SMH-TB, and the granularity of the data they contain. A detailed collection of
 160 all the variables available in the database is provided in Table 1.

161 **Fig 2: Patient-level and Encounter-level Data in SMH-TB**

162 **Table 1: Variables available in SMH-TB from both structured and unstructured sources**

Demographics	Tuberculosis Diagnosis
Patient ID	Known TB exposure*
MRN	BCG vaccination status*
Sex	TST performed*
Date of birth	TST induration*
Street address	TST interpretation*
Postal code ^a	IGRA performed*
Country of origin*	IGRA interpretation*
Year of immigration*	Diagnosis of active TB*
Immigration status	Diagnosis of LTBI*
Housing status	
Insurance status	Tuberculosis Medications
Patient is a healthcare worker*	Ever started isoniazid*
	Ever started rifampin*
Encounter Details	Ever started pyrazinamide*
Encounter ID	Ever started ethionamide*
Encounter type	Ever started vitamin B6*

Encounter date	
Direct cost ^b	Medical Conditions and Comorbidities
Indirect cost ^c	Autoimmune conditions ^{d*}
	Diabetes*
Aggregate Variables	Hematological malignancy*
Number of sputum inductions	Non-hematological malignancy*
Number of chest x-rays	Transplant performed*
Number of chest computed tomography	Renal failure ^{e*}
Hospital admission during course of TB outpatient care	Silicosis*
Number of emergency room visits during course of TB outpatient care	Hepatitis B
	Hepatitis C
Laboratory Results	HIV status*
AST	
ALT	Microbiology Reports**
CBC (Hb, Platelets, WBC)	Radiology Reports**
Cr	Pathology Reports**
Bilirubin	

163 MRN: Medical record number; AST: Aspartate transaminase; ALT: Alanine transaminase; CBC:
 164 Complete blood count; Hb: Hemoglobin; WBC: White blood cells; Cr: Creatinine; TB: Tuberculosis;
 165 BCG: Bacillus Calmette–Guérin; TST: Tuberculin sensitivity test; IGRA: Interferon gamma release
 166 assay; LTBI: Latent tuberculosis infection; HIV: Human immunodeficiency viruses
 167 ^aThe database only stores the Forward Sortation Area portion of the postal code of the patient’s residence.
 168 ^bDirect cost corresponds to health care services directly associated with the patient’s care including all
 169 nursing, allied health, diagnostic and therapeutic services, pharmaceutical and medical/surgical supplies
 170 for each visit.
 171 ^cIndirect cost corresponds to administrative and support services performed on behalf of all patients
 172 including information system and housekeeping overheads.

173 ^dAutoimmune conditions include: Sjogren’s syndrome, arthropathy, spondyloarthropathy, psoriatic
174 arthritis, rheumatoid arthritis, reactive arthritis, mixed connective tissue disease, connective tissue disease,
175 systemic lupus erythematosus, CREST syndrome, dermatomyositis, Wegener’s granulomatosis,
176 Goodpasture syndrome, vasculitis and psoriasis.

177 ^eRenal failure includes: nephropathy, renal insufficiency and glomerulonephritis.

178 *Variables collected from unstructured dictations and reports using natural language rulesets

179 **Unstructured text from which variables will be generated using natural language rulesets

180 ***Removing identifiable information***

181 There are two versions of SMH-TB. The full version includes indelible patient identifiers such as
182 a patient’s provincial health insurance (Ontario Health Insurance Plan) number; their SMH-
183 specific medical record number; all patient encounters whose encounter record is specific to a
184 given patient; laboratory test records whose encounter record is also specific to a given patient;
185 and all unstructured text data per encounter per patient. The patient identifiers allow for a fully
186 linked database, which can be updated and linked via future data extraction. The identifiable
187 unstructured data are also retained to support the development and testing of additional natural
188 language rulesets.

189 The de-identified version of SMH-TB is the version that will be primarily used for research
190 studies. It excludes the unstructured data and has been stripped of the following: hospital patient
191 ID, hospital encounter ID, address and day and month of date of birth. Each patient and
192 encounter is then re-coded with new unique IDs, and with the age in years on the date of the first
193 TB clinic encounter

194 ***Patient identification and validation***

195 The Decision Support Services (DSS) at SMH identified encounters which were coded as
196 services provided in the TB outpatient clinic to identify all TB patients. We then randomly
197 selected a list of 200 patients seen in the TB outpatient clinic (using clinic schedules with unique
198 patient identifiers stored separately from the EDW) to manually validate the codes used by DSS
199 to identify TB clinic outpatients, and validated that all (100%) identified patients were registered
200 in the TB clinic. To ensure high specificity of our identification of TB clinic patients, we
201 examined additional metadata (such as a mention of the TB clinic in the patient’s dictations) and
202 removed patients without matching metadata. SMH-TB therefore may include the rare patient
203 where the clinic visit codes in the EDW erroneously labelled a visit as a TB clinic visit, but this
204 estimate is expected to be <0.2% because of the additional metadata checks. The hospital unique
205 patient identifier for each individual was then cross referenced to lists of all individuals with
206 inpatient stays and ER visits to derive TB patient data from all sites of contact for TB care.

207 ***Data transformation (unstructured text to structured variables)***

208 Unstructured clinician dictations were used to create patient-level variables on demographics, TB
209 diagnosis, TB medications and comorbidities. The data for these variables were extracted using

210 rule-based information extraction tool CHARTextract (22). CHARTextract uses regular
211 expressions in order to perform pattern matching on text. Regular expressions have been used to
212 perform data extraction and even classification due to their high expressivity (17,23,24). These
213 capabilities come at the cost of a complex syntax, and thus rule creation typically involves the
214 expertise of a clinician who understands the subject matter and an interpreter who can express
215 the idea into regular expression syntax. We created a tiered rule system, where primary rules are
216 used to filter text at the sentence level using a scoring system and secondary rules can be used to
217 further enhance the weighting of the sentence. The tool applies the user-created rules to the data
218 and extracts the variables on-the-fly. The interface displays mismatches between the tool
219 prediction and the gold-standard label. Users can iterate on the rule creation process, allowing for
220 easy refinement and quick development of the rules. Fig 3 shows a component of a ruleset for
221 extracting diagnosis of active tuberculosis.

222 **Fig 3: Example of a component of a ruleset for extracting a variable (active TB diagnosis)**
223 **from unstructured text in clinical dictations (using CHARTextract)**

224 In order to create the rulesets used by CHARTextract, two clinicians (JB, SM) from the TB
225 outpatient clinic were consulted on dictation language and style. Clinicians (JB, SM, AA, and
226 KC) and a medical student (CW) manually labeled dictations for 200 patients from a subset of
227 the dataset to be used for validation. The set of 200 patients was selected from consecutive clinic
228 visits based on registered patient lists external to the EHR. This was done using the QuickLabel
229 tool which provides a user interface for streamlined labelling of specific variables, as well as the
230 option to label multiple variables simultaneously (25). Refinement of the natural language
231 rulesets was done by comparing the labels extracted by the rulesets via CHARTextract with the
232 manual labels. The refined rulesets are available as a real-time source as additional variables
233 from unstructured data (microbiology, radiology, and pathology reports) are generated (20).

234 ***Evaluation of data extraction***

235 To measure the performance of our rulesets and evaluate their generalizability to unseen data, we
236 calculated accuracy, recall, precision and F_1 scores. Recall (sensitivity) measures the ability of
237 the classifier to correctly distinguish true positive from false negative examples. Precision
238 (positive predictive value) measures the ability of the classifier to correctly distinguish true
239 positive from false positive examples. The F_1 score computes a harmonic mean of precision and
240 recall. Recall, precision and F_1 score were averaged across variable labels.

241 ***Binomial proportions estimated from extracted variables***

242 We used the refined rulesets to extract variables from the full dataset of patients with at least 1
243 dictation (N=3237). We converted “Yes/No/Not recorded” and “Positive/Negative/Unknown/Not
244 recorded” variables into binary 0-1 variables by assigning a value of 1 to patients with an
245 extracted value of “Yes” or “Positive”, and a value of 0 otherwise. We estimated the proportion
246 and 95% confidence intervals of patients for which the rulesets extracted “Yes” or “Positive” for

247 these variables using two methods: (1) logistic regression model without covariates, and (2) MC-
248 SIMEX model that accounts for the misclassification error in the extracted variables that was
249 calculated from the set of 200 manually abstracted patients (26). Briefly, for a binary random
250 variable Y , we estimate the probability $P(Y = 1)$ using a logistic regression model without
251 covariates, given by:

$$P(Y = 1) = h(\beta_0)$$

252 where h is the logistic function. Under the MC-SIMEX model, the binary random variable was
253 observed with misclassification errors, denoted by Y^* . We estimate the probability $P(Y^* = 1)$ as:

$$P(Y^* = 1) = h(\beta_0^*)$$

254 where β_0^* is defined as:

$$\beta_0^*(\lambda) = h^{-1}[\pi_{11}^\lambda h(\beta_0) + (1 - \pi_{00}^\lambda)(1 - h(\beta_0))]$$

255 π_{00} and π_{11} denote the specificity and sensitivity of Y^* , respectively, and λ is the
256 misclassification parameter. The final estimate for β_0^* is computed by a simulation-extrapolation
257 procedure described in (26).

258 Results

259 Population

260 A patient overview based on demographics is presented in Table 2. 3298 patients were included
261 in the database. The median age of the patients is 45 years, with an interquartile range of 34 to
262 58. There is a higher percentage of females than males in the cohort, around 57%. At least 79%
263 of the clinic patients were born outside of Canada, based on data extracted from patients'
264 dictations. The vast majority of patients were adequately housed, with publicly funded provincial
265 health care insurance (OHIP).

266 **Table 2: Demographics of the patients included in the SMH-TB database, 2011-2018.**

Variable	Value	Number of patients who attended at least 1 clinic visit (Total N=3298)	
		Count	Percentage
Age-group in years (median: 45, IQR: 34-58)	10-20	7	0.212
	20-30	422	12.8
	30-40	802	24.3

	40-50	705	21.4
	50-60	575	17.4
	60-70	388	11.8
	70-80	245	7.42
	80-90	126	3.82
	90-100	30	0.910
	100-110	2	0.0606
Sex	Female	1884	57.1
	Male	1417	42.9
	Missing ^a	1	0.0303
Born in Canada	Born in Canada	247	7.48
	Born outside Canada	2619	79.3
	Missing ^b	436	13.2
Underhoused ^c	Yes	80	2.42
	No	3222	97.6
Type of health insurance ^d	Ontario Health Insurance Plan (OHIP)	2859	86.6
	Uninsured Person Program (TB-UP)	221	6.69
	Refugee Health Coverage	78	2.36
	University Health Insurance Plan (UHIP)	41	1.24
	Self-payd	76	2.30
	Other ^e	27	0.819

267 IQR: Interquartile range

268 ^aMay be due to error in data entry at time of patient registration.

269 ^bPatient dictations did not mention immigration status or country of birth, or no dictations were found.

270 ^cUnderhoused: includes patients living in homeless shelters, group homes or patients with no fixed
 271 address.
 272 ^dFor patients with more than one type of insurance, only the insurance type used for the latest encounter is
 273 displayed in this table.
 274 ^eIncludes any patients with an out-of-province insurance, or not recorded insurance type.
 275

276 ***Evaluation of data extraction***

277 A summary of the rulesets' performance metrics for the 25 variables extracted from unstructured
 278 dictations is presented in Table 3. Diagnosis of active TB and LTBI rulesets had 97.5% and 96%
 279 accuracy, and 97.4% and 94.7% F₁ score, respectively. Rulesets for extracting TB medications
 280 generally achieved above 90% accuracy, recall and precision metrics.

281 **Table 3: Summary of performance metrics on test set for variables extracted from**
 282 **unstructured dictations. Patients included in test set: N = 200.**

Variable	True Positive*	True Negative*	Accuracy	Recall	Precision	F ₁ Score
Demographics						
Country of origin	--	--	0.970	0.987	0.987	0.986
Year of immigration	--	--	0.805	0.834	0.891	0.850
Patient is a healthcare worker	29	171	0.940	0.850	0.897	0.871
Tuberculosis Diagnosis						
Known TB exposure	43	157	0.965	0.952	0.945	0.949
BCG vaccination status	89	111	0.865	0.852	0.887	0.859
TST performed	100	100	0.990	0.990	0.990	0.990
TST induration	--	--	0.985	0.954	0.960	0.957
TST interpretation	86	114	0.980	0.978	0.981	0.980
IGRA performed	14	186	1.00	1.00	1.00	1.00

IGRA interpretation	5	195	1.00	1.00	1.00	1.00
Diagnosis of active TB	120	80	0.975	0.975	0.973	0.974
Diagnosis of LTBI	49	151	0.960	0.953	0.941	0.947
Tuberculosis Medications						
Ever started isoniazid	150	50	0.960	0.933	0.959	0.945
Ever started rifampin	127	73	0.970	0.962	0.974	0.967
Ever started pyrazinamide	124	76	0.995	0.996	0.994	0.995
Ever started ethambutol	118	82	0.985	0.985	0.984	0.984
Ever started vitamin B6	147	53	0.990	0.987	0.987	0.987
Medical Conditions and Comorbidities**						
Autoimmune conditions	8	192	0.965	0.862	0.767	0.807
Diabetes	26	174	0.945	0.870	0.883	0.876
Hematological malignancy	2	198	0.990	0.748	0.748	0.748
Non-hematological malignancy	12	188	0.955	0.937	0.787	0.843
Renal failure	8	192	0.975	0.807	0.849	0.827
HIV status	2	198	0.995	0.998	0.833	0.899

283 TB: Tuberculosis; BCG: Bacillus Calmette–Guérin; TST: Tuberculin sensitivity test; IGRA: Interferon
 284 gamma release assay; LTBI: Latent tuberculosis infection; HIV: Human immunodeficiency viruses

285 *True positives are defined as observations with a value of “Yes” or “Positive”; True negatives are
 286 defined as the complement of true positives; only applicable for extracted variables which have
 287 “Yes/No/Not recorded” or “Positive/Negative” values.
 288 **Patients that had undergone a transplant and patients diagnosed with silicosis were excluded from this
 289 table due to having no positive example in the test set.

290 ***Binomial proportions estimated from extracted variables***

291 The estimated proportions and their 95% confidence intervals created from the “Yes/No/Not
 292 recorded” and “Positive/Negative” extracted variables are given in Table 4.

293 **Table 4: Binomial proportion estimate and 95% confidence interval (CI) using standard**
 294 **binary regression and MC-SIMEX model for binary variables created from extracted**
 295 **variables. Total patients with at least 1 dictation: N = 3237.**

Description	Count (N=3237)	Logistic regression estimate (95% CI)	MC-SIMEX model estimate (95% CI)
Demographics			
Healthcare workers	438	13.5% (12.4, 14.8)	2.48% (2.02, 3.04)
Tuberculosis Diagnosis			
Known TB exposure	706	21.8% (20.4, 23.3)	16.8% (15.3, 18.3)
Received BCG vaccination	1316	40.7% (39.0, 42.4)	24.8% (23.0, 26.7)
Performed a TST	2279	70.4% (68.8, 72.0)	69.7% (68.1, 71.3)
Received a positive TST interpretation	2031	62.7% (61.1, 64.4)	61.9% (60.2, 63.6)
Performed an IGRA	296	9.14% (8.20, 10.2)	9.14% (8.20, 10.2)
Received a positive IGRA interpretation	301	9.30% (8.35, 10.3)	9.30% (8.35, 10.3)
Diagnosed with active TB	640	19.8% (18.4, 21.2)	18.2% (16.8, 19.7)
Diagnosed with LTBI	1473	45.5% (43.8, 47.2)	39.7% (37.8, 41.6)

Tuberculosis Medications			
Ever started on isoniazid	1314	40.6% (38.9, 42.3)	45.6% (43.6, 47.5)
Ever started on rifampin	548	16.9% (15.7, 18.3)	17.6% (16.3, 19.1)
Ever started on pyrazinamide	349	10.8% (9.76, 11.9)	9.99% (8.96, 11.1)
Ever started on ethambutol	348	10.8% (9.73, 11.9)	9.36% (8.32, 10.5)
Ever started on vitamin B6	986	30.5% (28.9, 32.1)	30.6% (29.0, 32.2)
Medical Conditions and Comorbidities*			
Autoimmune conditions	167	5.16% (4.45, 5.98)	0.259% (0.175, 0.383)
Diabetes	179	5.53% (4.79, 6.37)	0.358% (0.247, 0.517)
Hematological malignancy	71	2.19% (1.74, 2.76)	0.00625% (0.00320, 0.0122)
Non-hematological malignancy	140	4.32% (3.68, 5.08)	0.860% (0.599, 1.23)
Renal failure	65	2.01% (1.58, 2.55)	0.00895% (0.00450, 0.0180)
Diagnosed with HIV	175	5.41% (4.68, 6.24)	5.43% (4.69, 6.26)
No relevant medical conditions/comorbidities**	2569	79.4% (77.9, 80.7)	89.3% (87.9, 90.6)

296 MC-SIMEX: Misclassification Simulation Extraction; CI: Confidence interval; TB: Tuberculosis; BCG:
 297 Bacillus Calmette–Guérin; TST: Tuberculin sensitivity test; IGRA: Interferon gamma release assay;
 298 LTBI: Latent tuberculosis infection; HIV: Human immunodeficiency viruses
 299 *Patients that had undergone a transplant and patients diagnosed with silicosis were excluded from this
 300 table due to having no positive example in the test set.

301 **Includes any patient with an extracted value of “No/Not recorded/Negative” for all medical
302 conditions/comorbidities listed in the table.

303 After accounting for misclassification errors, the proportion of patients with an active TB
304 diagnosis was 18.2% and the proportion of patients with an LTBI diagnosis was 39.7%. 69.7%
305 of patients had performed a tuberculin sensitivity test and 61.9% of all patients had a positive
306 result for the test. The proportions of patients who were ever started on isoniazid, rifampin or B6
307 were 45.6%, 17.6% and 30.6% percent, respectively.

308 Discussion

309 To facilitate research on TB clinical epidemiology, diagnostics, clinical care and program
310 implementation, quality improvement, and linkage for future therapeutics trials and biomarker
311 studies, we developed a retrospective database of TB clinic patients using structured and
312 unstructured EHR data. The cohort and database are unique in the transformation of unstructured
313 data into structured variables using natural language rulesets with excellent performance when
314 validated against manual chart abstraction. The rulesets are open access, and the database is
315 accessible for research and open for collaboration with approval from local research ethics board.

316 The strength of the SMH-TB database comes from the inclusion of granular data, achieved by
317 extracting it from unstructured sources using natural language processing. While the database
318 contains standard structured data accessible in a wide variety of EHRs, a large and unique
319 component of our data comes directly from unstructured dictated clinic notes, which contain a
320 vast number of variables that can be used for a broad range of research topics, including, for
321 example, clinical epidemiology and modelling studies. The NLP rulesets allow us to create
322 granular patient-level variables from unstructured data accurately and efficiently, reducing the
323 amount of time spent on manual abstraction to a minimum. Moreover, the large amount of
324 unstructured raw data is a tremendous resource for evaluating and deploying machine learning
325 and deep learning models capable of automatically extracting meaningful variables from clinical
326 notes (27–29). Machine learning models and workflows can be developed to leverage the
327 structured and extracted variables for predictive modeling and early warning systems (30–32).
328 The breadth of data provided makes this a unique and powerful tool in both clinical and
329 computational research.

330 The main limitations of the SMH-TB database include issues that arise from missing or incorrect
331 data and the limited availability of data for certain variables leading to non-robust natural
332 language rulesets. Data errors can be due to both human and algorithmic mistakes. Much of the
333 burden of including relevant data in clinical dictations lies with the clinician attending the patient
334 and dictating the note. In the absence of a standardized format, as was the case in the SMH TB
335 clinic, variables may not be dictated in a manner that enables their capture by the NLP tools, or
336 are not dictated at all. The creation of a shared set of guidelines and standard formatting for TB

337 clinic dictations, containing all variables relevant to the database, will ensure all data required are
338 captured with future database updates.

339 When the unstructured data undergoes information extraction, mislabeling of variables can occur
340 due to certain rulesets having subpar performance. This issue is especially apparent for variables
341 with scarce availability of labels. For example, in our validation dataset there were no patients
342 with silicosis. The ruleset for classifying silicosis was adapted from other immunosuppressive
343 conditions and expert knowledge in disease. While it is possible that such rulesets are overly
344 confident in assigning a “No” label to patients even if they present with the condition in
345 question, given the rarity of the event in the patient population it was not possible to provide
346 further cases for perfection of refinement of the NLP ruleset. As such, we have indicated the
347 metrics of our variables (Table 3), so that researchers can understand the limitations of the data
348 with which they are working. The 200 charts sampled for ruleset refinement were consecutive
349 patients from a set of clinic visits and may not have been sufficient for less common variables
350 such as comorbidities. That is, further ruleset refinement will be needed with additional charts
351 with purposive sampling of true positives of infrequent variables.

352 **Conclusion**

353 In summary, here we describe the SMH-TB cohort and database which aim to be a resource for
354 scientists who are conducting research into many facets of TB. The database is unique in that it
355 contains highly granular socio-demographic and clinical patient data derived from structured and
356 unstructured EHR data extracted using NLP rule sets. The validated rule sets are provided open
357 access for use and the data base is intended to be available for collaborative studies.

358 **Data Availability**

359 The validated NLP rulesets are publicly available for use from: [https://github.com/mishra-lab/tb-](https://github.com/mishra-lab/tb-nlp-rulesets)
360 [nlp-rulesets](https://github.com/mishra-lab/tb-nlp-rulesets). Data collected in SMH-TB contains sensitive patient information and as such,
361 researchers interested in conducting TB-related research using the data are welcome to contact
362 the corresponding author and submit a request. The study team welcomes collaboration and use
363 of the database, and all external requests will be screened to ensure adequate data exists to enable
364 a collaboration. The project will then undergo the approval process of the Research and Ethics
365 Board (REB) of Unity Health Toronto. Data provided to researchers can either be the de-
366 identified version of the SMH-TB database, or the full identifiable version, based on their
367 research needs and REB approval.

368 **Funding**

369 Supported by the Ontario Early Researcher Award Number ER17-13-043 (to SM). The funders
370 had no role in study design, data collection and analysis, decision to publish, or preparation of
371 the manuscript.

372 **Acknowledgements**

373 We thank Dr. Natasha Sabur for supporting arbitration for rulesets; Julie Seemangal (TB
374 Outpatient Clinic Co-Lead) and Grace Bezaliel for supporting verification of algorithms to
375 classify patients seen in the TB clinic.

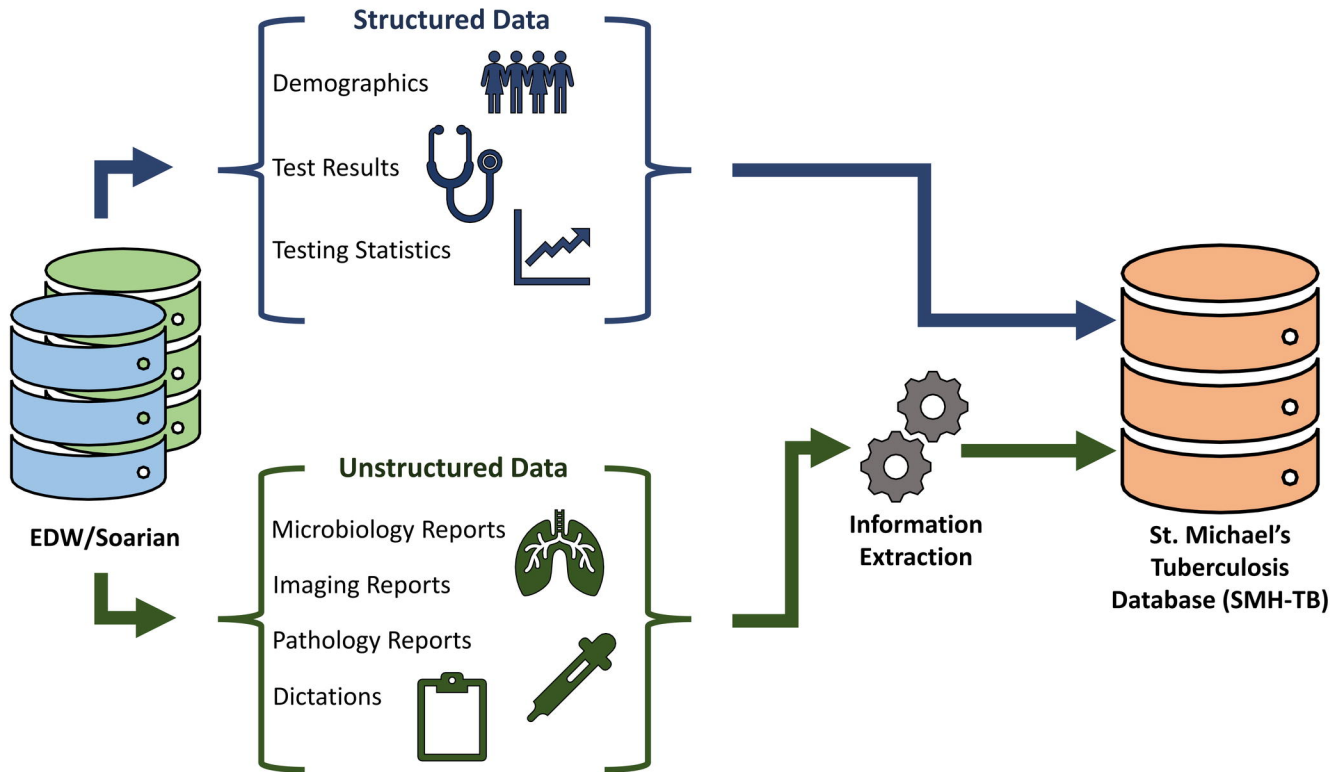
376

377 **References**

- 378 1. Reid MJA, Arinaminpathy N, Bloom A, Bloom BR, Boehme C, Chaisson R, et al. Building a
379 tuberculosis-free world: The Lancet Commission on tuberculosis. *Lancet Lond Engl*. 2019
380 Mar 30;393(10178):1331–84.
- 381 2. Uplekar M, Weil D, Lönnroth K, Jaramillo E, Lienhardt C, Dias HM, et al. WHO’s new End
382 TB Strategy. *The Lancet*. 2015 May 2;385(9979):1799–801.
- 383 3. Lönnroth K, Migliori GB, Abubakar I, D’Ambrosio L, Vries G de, Diel R, et al. Towards
384 tuberculosis elimination: an action framework for low-incidence countries. *Eur Respir J*.
385 2015 Apr 1;45(4):928–52.
- 386 4. CDC. Deciding When to Treat Latent TB Infection [Internet]. 2018 [cited 2020 Aug 25].
387 Available from: <https://www.cdc.gov/tb/topic/treatment/decideltbi.htm>
- 388 5. Kim PS, Makhene M, Sizemore C, Hafner R. Viewpoint: Challenges and Opportunities in
389 Tuberculosis Research. *J Infect Dis*. 2012 May 15;205(suppl_2):S347–52.
- 390 6. Busingye D, Gianacas C, Pollack A, Chidwick K, Merrifield A, Norman S, et al. Data
391 Resource Profile: MedicineInsight, an Australian national primary health care database. *Int J*
392 *Epidemiol*. 2019 Dec 1;48(6):1741–1741h.
- 393 7. Garies S, Birtwhistle R, Drummond N, Queenan J, Williamson T. Data Resource Profile:
394 National electronic medical record data from the Canadian Primary Care Sentinel
395 Surveillance Network (CPCSSN). *Int J Epidemiol*. 2017 Aug 1;46(4):1091–1092f.
- 396 8. Finer S, Martin HC, Khan A, Hunt KA, MacLaughlin B, Ahmed Z, et al. Cohort profile: East
397 London genes & health (ELGH), a community-based population genomics and health study
398 of British Bangladeshi and British Pakistani people. *Int J Epidemiol* [Internet]. [cited 2020
399 Mar 2]; Available from: [https://academic.oup.com/ije/advance-](https://academic.oup.com/ije/advance-article/doi/10.1093/ije/dyz174/5555939)
400 [article/doi/10.1093/ije/dyz174/5555939](https://academic.oup.com/ije/advance-article/doi/10.1093/ije/dyz174/5555939)
- 401 9. Ashfaq A, Lönn S, Nilsson H, Eriksson JA, Kwatra J, Yasin ZM, et al. Data resource profile:
402 Regional healthcare information platform in Halland, Sweden, a dedicated environment for
403 healthcare research. *Int J Epidemiol* [Internet]. [cited 2020 Mar 2]; Available from:
404 <https://academic.oup.com/ije/advance-article/doi/10.1093/ije/dyz262/5701527>
- 405 10. Office of the National Coordinator for Health Information Technology. Office-based
406 Physician Electronic Health Record Adoption [Internet]. 2019 [cited 2020 Apr 7]. Available
407 from: dashboard.healthit.gov/quickstats/pages/physician-ehr-adoption-trends.php
- 408 11. Henry J, Pylypchuk Y, Searcy T, Patel V. Adoption of electronic health record systems
409 among US non-federal acute care hospitals: 2008–2015. *ONC Data Brief*. 2016;35:1–9.

- 410 12. Chan KS, Fowles JB, Weiner JP. Review: Electronic Health Records and the Reliability and
411 Validity of Quality Measures: A Review of the Literature. *Med Care Res Rev*. 2010 Oct
412 1;67(5):503–27.
- 413 13. Nicholson A, Tate AR, Koeling R, Cassell JA. What does validation of cases in electronic
414 record databases mean? The potential contribution of free text. *Pharmacoepidemiol Drug*
415 *Saf*. 2011;20(3):321–4.
- 416 14. Khan K, Campbell A, Wallington T, Gardam M. The impact of physician training and
417 experience on the survival of patients with active tuberculosis. *CMAJ Can Med Assoc J*.
418 2006 Sep 26;175(7):749–53.
- 419 15. Long R, Heffernan C, Gao Z, Egedahl ML, Talbot J. Do “Virtual” and “Outpatient” Public
420 Health Tuberculosis Clinics Perform Equally Well? A Program-Wide Evaluation in Alberta,
421 Canada. *PLoS ONE* [Internet]. 2015 Dec 23 [cited 2020 Apr 7];10(12). Available from:
422 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4689372/>
- 423 16. Baldwin KB. Evaluating Healthcare Quality Using Natural Language Processing. *J Healthc*
424 *Qual*. 2008;30(4):24–9.
- 425 17. Luo Y, Thompson WK, Herr TM, Zeng Z, Berendsen MA, Jonnalagadda SR, et al. Natural
426 Language Processing for EHR-Based Pharmacovigilance: A Structured Review. *Drug Saf*.
427 2017;40(11):1075–89.
- 428 18. Wang Y, Wang L, Rastegar-Mojarad M, Moon S, Shen F, Afzal N, et al. Clinical
429 information extraction applications: A literature review. *J Biomed Inform*. 2018 Jan 1;77:34–
430 49.
- 431 19. WHO. Tuberculosis country profiles [Internet]. World Health Organization; [cited 2020 Apr
432 7]. Available from: <http://www.who.int/tb/country/data/profiles/en/>
- 433 20. Landsman D, LKS-CHART. Tuberculosis NLP Rulesets [Internet]. GitHub. [cited 2020 Jul
434 1]. Available from: <https://github.com/mishra-lab/tb-nlp-rulesets>
- 435 21. Government of Canada. Tri-Council Policy Statement: Ethical Conduct for Research
436 Involving Humans – TCPS 2 (2018) [Internet]. 2019 [cited 2020 Aug 27]. Available from:
437 https://ethics.gc.ca/eng/policy-politique_tcps2-eptc2_2018.html
- 438 22. LKS-CHART. CHARTextract [Internet]. CHARTextract. [cited 2020 Jun 29]. Available
439 from: <https://lks-chart.github.io/CHARTextract-docs/>
- 440 23. Rosier A, Burgun A, Mabo P. Using regular expressions to extract information on pacemaker
441 implantation procedures from clinical reports. *AMIA Annu Symp Proc*. 2008;2008:81–5.
- 442 24. Finley G, Edwards E, Robinson A, Brenndorfer M, Sadoughi N, Fone J, et al. An automated
443 medical scribe for documenting clinical encounters. In: *Proceedings of the 2018 Conference*
444 *of the North American Chapter of the Association for Computational Linguistics:*
445 *Demonstrations* [Internet]. New Orleans, Louisiana: Association for Computational

- 446 Linguistics; 2018 [cited 2020 Jun 29]. p. 11–15. Available from:
447 <https://www.aclweb.org/anthology/N18-5003>
- 448 25. Joshi U. QuickLabel [Internet]. 2019 [cited 2020 Jun 29]. Available from:
449 <https://github.com/Sabrewarrior/QuickLabel>
- 450 26. Küchenhoff H, Mwalili SM, Lesaffre E. A General Method for Dealing with
451 Misclassification in Regression: The Misclassification SIMEX. *Biometrics*. 2006;62(1):85–
452 96.
- 453 27. Jagannatha AN, Yu H. Structured prediction models for RNN based sequence labeling in
454 clinical text. *Proc Conf Empir Methods Nat Lang Process Conf Empir Methods Nat Lang*
455 *Process*. 2016 Nov;2016:856–65.
- 456 28. Wu Y, Jiang M, Lei J, Xu H. Named Entity Recognition in Chinese Clinical Text Using
457 Deep Neural Network. *Stud Health Technol Inform*. 2015;216:624–8.
- 458 29. Fries J. Brundefly at SemEval-2016 Task 12: Recurrent Neural Networks vs. Joint Inference
459 for Clinical Temporal Information Extraction. In: *Proceedings of the 10th International*
460 *Workshop on Semantic Evaluation (SemEval-2016)* [Internet]. San Diego, California:
461 Association for Computational Linguistics; 2016 [cited 2020 Aug 25]. p. 1274–1279.
462 Available from: <https://www.aclweb.org/anthology/S16-1198>
- 463 30. Miotto R, Li L, Kidd BA, Dudley JT. Deep Patient: An Unsupervised Representation to
464 Predict the Future of Patients from the Electronic Health Records. *Sci Rep*. 2016 May
465 17;6(1):26094.
- 466 31. Tran T, Nguyen TD, Phung D, Venkatesh S. Learning vector representation of medical
467 objects via EMR-driven nonnegative restricted Boltzmann machines (eNRBM). *J Biomed*
468 *Inform*. 2015 Apr 1;54:96–105.
- 469 32. Choi E, Bahadori MT, Schuetz A, Stewart WF, Sun J. Doctor AI: Predicting Clinical Events
470 via Recurrent Neural Networks. *JMLR Workshop Conf Proc*. 2016 Aug;56:301–18.
- 471



Patient-Level

Encounter-Level

Diagnoses

Outpatient Visits

Procedures

Inpatient Visits

Raw Dictations

Emergency Visits

Lab Results

Imaging Results

Extracted Variables

Demographics

TB Diagnosis

TB Medications

Medical Conditions and
Comorbidities

No Yes +

Delete File Save File

Add Primary Rule

1

If active × appears, then score 0 + ×

and mycobacterium bovis × OR × tuberculosis × OR × \bTB\b × OR × infection × OR × disease × appears after 1, then score 3 ×

unless low (suspicion|chance)? × OR × rul(ing|es|ed|e)? out × OR × denies × OR × asymptomatic × appears, then score ×

unless past × OR × previous × OR × prior × appears, then score ×