

**Supplementary File.**

**Data Quality Assessment Documentation and Reporting Recommendations** (as outlined by Khan [25])

Data Capture		
Original Data Source		
Data origin	<p><b>NAPCRen:</b> EMR data from primary care clinical practices in northern Alberta using one of Accuro (<i>n</i>=1), Healthquest (<i>n</i>=1), Med Access (<i>n</i>=11), or Wolf (<i>n</i>=6) EMR system.</p> <p>As of December 31, 2018, data was captured from: 19 primary care clinics 79 providers (family physicians) 104,622 patients</p>	<p><b>SAPCRen:</b> EMR data from primary care clinical practices in southern Alberta using one of Med Access (<i>n</i>=12), Practice Solutions (<i>n</i>=1), or Wolf (<i>n</i>=22) EMR system.</p> <p>As of December 31, 2018, data was captured from: 34 primary care clinics + 1 community pediatric clinic 239 providers (family physicians, nurse practitioners, community pediatricians) 306,329 patients</p>
Data capture method	<p>Direct entry by clinicians/clinical staff into EMR system; interface and available templates/data fields may vary by EMR system.</p> <p>Automated upload of laboratory results into EMR from community lab provider.</p>	
Original collection purpose	<p>Patient clinical care and administrative tasks (i.e. billing) in community-based primary care clinics in Alberta.</p>	
Data Steward Information		
Data steward	<p><b>NAPCRen:</b> Northern Alberta Primary Care Research Network (NAPCRen), one of 11 regional practice-based research networks hosting the CPCSSN project. Housed within the Department of Family Medicine, University of Alberta; Edmonton, Alberta, Canada.</p>	<p><b>SAPCRen:</b> Southern Alberta Primary Care Research Network (SAPCRen), one of 11 regional practice-based research networks hosting the CPCSSN project. Housed within the Department of Family Medicine, University of Calgary; Calgary, Alberta, Canada.</p>

Database model/structure	<p>The CPCSSN data are stored as a SQL database (either SQL Server or SQLite format) and currently contains 26 tables:</p> <ul style="list-style-type: none"> <li>▪ Eight tables house clinic, health care provider, or administrative information (e.g. patient-provider assignment, non-identifiable clinic information, database version information)</li> <li>▪ Fifteen tables hold both original EMR data and processed (cleaned and coded) data</li> <li>▪ Two tables (DiseaseCase, DiseaseCaseIndicator) contain a list of patients who have been indexed with one or more of the conditions for which CPCSSN has a validated definition; these definitions use data derived from the 15 data tables</li> <li>▪ The Deprivation table is empty, but will be populated in the future.</li> </ul> <p>Each table is assigned an integer primary key labelled: &lt;table_name&gt;_ID. Most data tables also contain the columns:</p> <ul style="list-style-type: none"> <li>▪ Site_ID (an integer assigned to each primary care practice) and Network_ID (an integer assigned to each PBRN) together form a composite primary key in the Site table, and a composite foreign key in most CPCSSN tables</li> <li>▪ Patient_ID (randomly assigned integer for each patient) is the primary key in Patient table and a foreign key in most other tables</li> <li>▪ Encounter_ID is the integer primary key for the Encounter table, but a foreign key in several other tables</li> <li>▪ Provider_ID is the integer primary key for the Provider table and a foreign key for several tables</li> <li>▪ All columns are strictly typed, enforced either by the database schema (SQL Server) or external software (SQLite)</li> </ul> <p>The Entity Relationship Diagram (ERD), which outlines relationships between all CPCSSN tables, can be found on the CPCSSN website: <a href="http://cpcssn.ca/research-resources/cpcssn-data-dictionary-and-erd/">http://cpcssn.ca/research-resources/cpcssn-data-dictionary-and-erd/</a></p> <p>For SAPCReN and NAPCReN, extraction, transformation, coding, and cleaning algorithms are programmed in Python.</p>
Data dictionary	<p>The most recent data dictionary is posted on the CPCSSN website: <a href="http://cpcssn.ca/research-resources/cpcssn-data-dictionary-and-erd/">http://cpcssn.ca/research-resources/cpcssn-data-dictionary-and-erd/</a></p>
<b>Data Processing / Data Provenance</b>	
Data extraction specifications	<p>The general procedure for EMR data extraction in both SAPCReN and NAPCReN is as follows:</p> <ul style="list-style-type: none"> <li>▪ Access the EMR frontend or backend, either through the vendor or by a CPCSSN data manager <ul style="list-style-type: none"> <li>○ <b>Wolf &amp; PS Suite:</b> backend extraction completed by Telus Health</li> <li>○ <b>Accuro &amp; Healthquest:</b> backend extraction conducted by a CPCSSN data manager using a log-in account for both clinic and EMR database</li> <li>○ <b>Med Access:</b> frontend extraction by a CPCSSN data manager using a log-in account created by the clinic</li> </ul> </li> <li>▪ Select data tables and columns of interest, then: <ul style="list-style-type: none"> <li>○ Null any columns with identifiable information (e.g. names, addresses, health care numbers)</li> <li>○ Select patients with an assigned CPCSSN provider or for those who have not been assigned any provider, use the four-cut method [31] to assign a provider</li> <li>○ An EMR mapping file is generated which consists of four data items: EMR patient ID (number assigned by each clinic EMR system), Network ID (number that uniquely identifies each CPCSSN regional network), Site ID (number that uniquely identifies each clinic within a network) and CPCSSN patient ID (unique patient number randomly assigned by CPCSSN). This file is used to ensure CPCSSN patient IDs are assigned uniquely to each EMR patient ID within each clinic and is held separately from the patient health data. The file can also be used to assist with re-identification for practice quality improvement or linkage with other data sources</li> </ul> </li> <li>▪ Transfer data to CPCSSN computers via a secure transfer method (e.g. SFTP)</li> </ul>

	<p>See Figure 2 for CPCSSN data pipeline.</p>
<p>Mappings from original values to standardized values</p>	<p>The CPCSSN database includes both original data extracted from the EMR, and cleaned or coded information resulting from CPCSSN’s algorithms (see CPCSSN Data Dictionary [28]). All original data are extracted from the EMR unchanged, except for a small number of EMR-specific code conversions and some parsing.</p> <p><u>General cleaning steps</u></p> <ul style="list-style-type: none"> <li>▪ Null empty strings and remove leading or trailing whitespace from all records</li> <li>▪ Normalize postal codes to format A1A 1A1 <ul style="list-style-type: none"> <li>○ NAPCReN: truncate to forward sortation area code (A1A)</li> </ul> </li> <li>▪ Convert all dates to format yyyy-mm-dd</li> <li>▪ Verify code formats for ATC, ICD-9, and DIN classification systems</li> <li>▪ Delete invalid labs where <ul style="list-style-type: none"> <li>○ lab name is a purely numeric value, or matches a text patterns that indicates it is not a useful lab name (e.g. ‘unit’, ‘patient name’, ‘date’)</li> <li>○ lab value non-numeric; exceptions are lab values starting with a comparison operator and Hepatitis C</li> </ul> </li> <li>▪ Convert medication dispensed forms and duration units to CPCSSN standard forms</li> <li>▪ Delete empty or duplicate records</li> <li>▪ Remove orphaned records (e.g. records with no matching Patient_ID in the Patient table)</li> <li>▪ Remove patients who only have demographic information and no other data</li> <li>▪ Parse billing records that contain multiple conditions into individual records</li> </ul> <p><u>Allowable ranges of values for NAPCReN and SAPCReN:</u></p> <ul style="list-style-type: none"> <li>▪ SBP: 50 to 300 mmHg</li> <li>▪ DBP: 20 to 200 mmHg</li> <li>▪ Weight: 1 kg to 500 kg</li> <li>▪ Height: 30 cm to 230 cm</li> <li>▪ BMI (adult): 5 to 200 kg/m<sup>2</sup></li> <li>▪ Waist circumference: 10 to 300 cm</li> <li>▪ Waist-to-hip ratio: 0.1 to 1</li> <li>▪ PEFr: 0-200 L/min</li> <li>▪ Ranges for laboratory values are specific to each of the 44 lab results processed by CPCSSN</li> <li>▪ Age at onset: ≤120</li> <li>▪ Birth year: extraction year (e.g. 2018) minus 120 years</li> <li>▪ Record creation date: earliest EMR date (1990-01-01) to extraction end date (2018-12-31)</li> <li>▪ Condition onset or procedure performed date: 1898-12-31 (extraction end date minus 120 years) to 2018-12-31</li> </ul> <p><u>Coding to reference system</u></p> <ul style="list-style-type: none"> <li>▪ <b>Medication:</b> existing code, DIN, or free text mapped to codes in the ATC Classification System</li> <li>▪ <b>Laboratory values:</b> existing code or free text mapped to LOINC; free-text mapping is completed for 44 labs</li> <li>▪ <b>Diagnostic codes:</b> ICD-9 or other codes (e.g. ICD-10, ICPC) found in Health Condition/Profile, Encounter Diagnosis, and Billing tables mapped to ICD-9 codes/classification heading</li> <li>▪ <b>Diagnostic text:</b> Pattern matches are used to map free text to ICD-9 codes; this mapping can be completed for 48 conditions in the Health Condition/Profile, Encounter Diagnosis, Billing tables</li> <li>▪ <b>Referrals:</b> original text in Referral field mapped to SNOMED concept code; not extracted from referral letters or any PDF documents</li> <li>▪ <b>Risk factors:</b> information in the Risk Factor table categorized by text pattern matching to one of the following headings: Smoking, Alcohol, Exercise, Diet, Obesity, or Psychosocial Stress.</li> </ul>

	<p>Smoking and Alcohol information further parsed (when available) into Status, Dose, Frequency, Duration, End date</p> <ul style="list-style-type: none"> <li>▪ <b>Family history:</b> relationship type is identified by pattern matching as one of: Mother, Father, Daughter, Son, (Half) Sister, (Half) Brother, (Great) Grandmother, (Great) Grandfather, (Great) Aunt, (Great) Uncle, Niece, Nephew, Cousin, or None. Relationship side (Maternal or Paternal) is also identified by text pattern matching</li> </ul>
<p>Data management organization's data transformation routines, including constructed variables</p>	<p>CPCSSN Data Dictionary [28] highlights CPCSSN-constructed variables in blue, including:</p> <ul style="list-style-type: none"> <li>▪ <b>Calculated age:</b> current year minus year of birth (<i>only as part of the national data provided to approved data users</i>)</li> <li>▪ <b>BMI:</b> uses original BMI if available; otherwise, CPCSSN algorithms generate BMI from original height and weight data where none previously exists</li> <li>▪ <b>Deprivation category:</b> maps dissemination areas (geographical units linked to postal code) to material and social quintiles [29] (<i>for future use in late 2019</i>)</li> <li>▪ <b>Disease case and disease case indicator tables:</b> CPCSSN has developed and validated case definitions based on combinations of diagnoses, billing codes, medications and lab results; patients are indexed according to these definitions. Definitions for conditions currently include: diabetes, hypertension, depression, osteoarthritis, chronic obstructive pulmonary disease, dementia, parkinsonism, epilepsy, plus chronic kidney disease, herpes zoster, and pediatric asthma [5, 6, 14]</li> </ul> <p>CPCSSN redacts several types of text that could potentially identify patients, providers, or clinics.</p> <ul style="list-style-type: none"> <li>▪ Non-medical words that may identify specific CPCSSN patients and providers are compared against exclusion lists, containing common words and/or allergy, diagnosis, and/or medication-specific terms to determine which are to be redacted. The exclusion list(s) to apply depend on the type of data being deidentified. Only those names not in the exclusion list(s) will be redacted in free text.</li> <li>▪ Other identifiers, such as social insurance numbers, province health numbers, driver's license numbers, credit card numbers, Worker's Compensation Board numbers, email addresses, and phone numbers (international, North American, and local) are suppressed</li> </ul>
<p>Data processing validation routines</p>	<p>During <b>extractions</b>, the CPCSSN data are verified by:</p> <ul style="list-style-type: none"> <li>▪ Comparing the number of files and data volume between extraction cycles occurring every six months</li> <li>▪ Patient and provider counts</li> <li>▪ Table row counts</li> <li>▪ Counts for CPCSSN-specific labs and exams</li> </ul> <p>During the <b>data transformation</b>, a time-stamped log-file provides an itemized description of all processing steps completed, including:</p> <ul style="list-style-type: none"> <li>▪ Initial and final row counts for all tables</li> <li>▪ Type(s) of transformation(s) and number of data elements transformed</li> <li>▪ Warnings and error messages (e.g. tables missing, processes that cannot be completed)</li> </ul> <p>During the <b>coding and cleaning processes</b>, the following validation tests are performed:</p> <ul style="list-style-type: none"> <li>▪ Time-stamped log-file which includes an itemized list of completed processes and warnings where data should be coded or is outside allowable ranges</li> <li>▪ Coder process verification tests used to compare the effect of software changes on specific coding or cleaning processes (e.g. ATC coding, exam data cleaning, deidentification) pre- and post-change</li> <li>▪ End-to-end test for case definition integration that compares changes in the validity metrics (i.e. sensitivity, specificity) for CPCSSN case definitions after any changes to any part of the coding and cleaning software</li> </ul>

Audit trail	<p>Separate time- and date-stamped log files are retained for each of the extraction, transformation, and cleaning/coding stages.</p> <ul style="list-style-type: none"><li>▪ <b>Extraction:</b> one file per clinic; contains lists of tables and columns extracted and any warnings or errors that occurred (e.g. connection lost, table missing). An introspection table is created detailing the properties of the source database (e.g. table/column names, data types).</li><li>▪ <b>Transformation:</b> one file per clinic; contains software version information, file locations, initial and final row counts for all tables, type(s) of transformation(s) and number of data elements transformed, warnings and error messages (e.g. tables missing, processes that cannot be completed)</li><li>▪ <b>Coding/Cleaning:</b> one per PBRN; includes an itemized list of processes completed, warnings where data cannot be coded or is outside of allowed ranges</li></ul>
-------------	---