# Rapid Universal Early Screening for Alzheimer's Disease and Related Dementia via Pattern Discovery in Diagnostic History

Dmytro Onishchenko[1], Sam Searle[7], Kenneth Rockwood[7], James A. Mastrianni[5,6] and Ishanu Chattopadhyay[1,2,3,4]★

[1]Department of Medicine, University of Chicago, Chicago, IL USA
[2]Committee on Genetics, Genomics & Systems Biology, University of Chicago, Chicago, IL USA
[3]Committee on Quantitative Methods in Social, Behavioral, and Health Sciences, University of Chicago, Chicago, IL USA
[4]Center for Health Statistics, Department of Medicine, University of Chicago, Chicago, IL USA
[5]Department of Neurology, University of Chicago, Chicago, IL USA
[6]Committee on Neurobiology, University of Chicago, Chicago, IL USA
[7]Division of Geriatric Medicine, Department of Medicine, Department of Community Health and Epidemiology, School of Health Administration, Halifax, NS Canada

★To whom correspondence should be addressed: e-mail: ishanu@u chicago.edu.

## SUMMARY

**Alzheimer's disease (AD) is a progressive, incurable and ultimately fatal neurodegenerative condition. In this study, we introduce the Zero-burden Co-morbid Risk (ZCoR) score to screen for the future risk of AD and related dementia (ADRD) $1-10$ years before a clinical diagnosis. Requiring no new bloodwork or cognitive tests, ZCoR leverages uncharted comorbidity patterns, to potentially enable near-instantenous universal point-of-care screening of entire patient populations. In validation, ZCoR ($n = 729,018$) achieves out-of-sample AUC $> 90\%$ for predicting a diagnosis immediately after screening, an AUC $> 87\%$ for a diagnosis made one year earlier than in current practice, and maintaining over $> 80\%$ AUC for predictions made a decade earlier, irrespective of sex. We achieve high predictability in patients lacking any of the currently suspected risk factors; demonstrating effectiveness in cohorts at higher risk of missed diagnoses. Additionally, ZCoR can target mild cognitive impairment (MCI) with performance at par with questionnaire-based assessments (AUC $88-90\%$), maintaining high effectiveness (AUC $\approx 80\%$) for predicting impairment upto $3$ years into the future. Powered by stochastic learning algorithms that enhance standard machine learning, ZCoR enables discovery in electronic heath record databases, can reduce ADRD and MCI diagnostic delays, and the impact of socio-economic and demographic variables, with immediate impact on patient outcomes.**

## INTRODUCTION

**D**EMENTIA is an acquired loss of cognition in one or more domains including learning and memory, social cognition, language, executive function, complex attention, and perceptual motor function, severe enough to significantly diminish social or occupational function[1]. Affecting approximately 47 million people worldwide[1], including over 5.5 million in the United States[2,3]; and projected to be over 81 million worldwide by 2040[4], there is an immediate need to find effective interventions, and screening tools that enable them.

The most common cause of dementia, contributing to 60%–80% of cases, is believed to be Alzheimer's disease (AD), a progressive, fatal, and currently incurable neurodegenerative condition[5]. Based on patient years lived with disability plus years lost to premature mortality, AD ranked as the 6[th] most burdensome disease or injury in the US in 2016, up from 12[th] in 1990[6], and was implicated in over 250,000 deaths in 2018[5].

AD-related neuropathology appears to progress over years or decades, independently of the clinical course, suggesting lengthy asymptomatic, subclinical, and/or subtly symptomatic periods before a clinical diagnosis[7].

a. Prevalence in Truven Dataset: Male

b. Prevalence in Truven Dataset: Female

c. Cumulative probability of ADRD diagnosis in Truven dataset

Fig. 1: **Descriptive statistics of database.** Panels a and b show geographic locations of study participants. Panel c illustrates the distribution of patient ages at documented diagnosis, showing that risk starts increasing from around 60 years, matching known ADRD onset age characteristics[3]. Also, panel c illustrates that the empirical risk per patient is higher for males in the Truven dataset, although there are more females with ADRD in total (See Tables I and II), and more female patients in general in the database in the relevant age groups. The age-stratified prevalence in the Truven dataset align closely with prevalence numbers reported for the US in 2020[20].

Accurate screening for both current and future cases on the Alzheimer's clinical spectrum may be expected to lead to earlier detection of AD biomarkers, neuropathology and incipient cognitive impairment or dementia. In turn, acclerating diagnosis may provide several important benefits for patients, caregivers, healthcare providers, and society[1-3,8-11]: first, pharmacologic and non-pharmacologic interventions may be applied to slow progression of cognitive impairment, while cognition is relatively preserved. Second, use of tailored education strate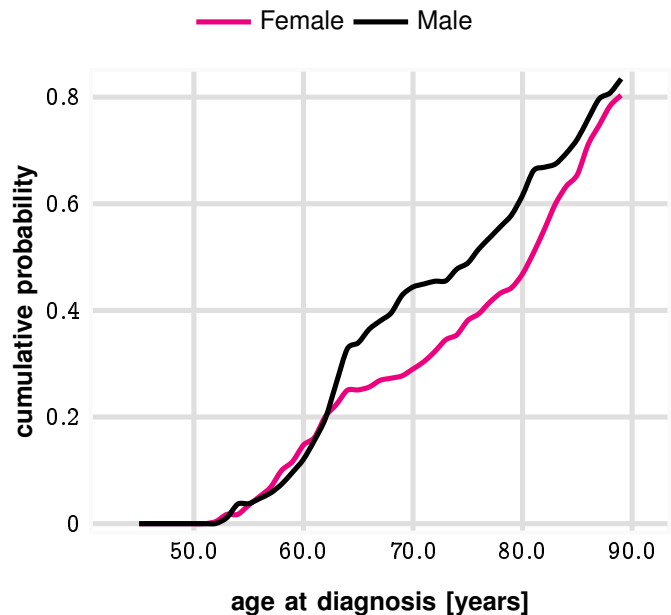gies may be facilitated to promote cognitively-impaired patients' adherence to and safe use of complex treatment regimens. Third, patient capacity for financial, legal, and health care decision-making may be optimized when it can occur as early as possible in the syndrome's course[9]. And finally, patient access to clinical trials of drugs targeting cognition and dementia may be fostered, and study samples may be enriched, accelerating progress and decreasing costs of such investigations.

However, accurate screening for ADRD is limited by the current diagnostic/prognostic modalities. Imaging or cerebrospinal fluid testing for evidence of beta-amyloid plaques and neurofibrillary tau deposits, the basis for an AD classification[7] and predictors of potential cognitive worsening, is expensive, invasive, and sometimes inaccessible. Although measurement of phosphorylated tau in plasma has shown promise as a specific marker and prognostic factor for ADRD[12-16], this method is as yet unavailable in everyday practice, entails an invasive blood draw, and if widely used, may in aggregate prove costly. Neuropsychological testing instruments such as the Montreal Cognitive Assessment (MOCA)[17,18] have good diagnostic accuracy and some prognostic utility in identifying mild cognitive impairment (MCI) and mild AD[19], but their time requirements, even when measured in minutes, may add appreciably to length-of-visit, and hence pose challenges in primary care settings[9,10]. Moreover, these instruments require validation when used in additional locales or languages, and efficacy in predicting future diagnoses might be suspect.

Analysis of routinely-collected health care data in past medical encounters may offer a passive, non-invasive,

TABLE I: Inclusion/Exclusion, Positive/Control Criteria & Cohort Definitions

| | **Definitions** |
|---|---|
| Inclusion/Exclusion Criteria | Age 50+ years |
| | Has medical history for $\geq$ 3 yrs |
| Cohort Definition | **Positive Cohort:** Patients either with at least one target code for ADRD from Tab. III (Case Dx), or with at least one of the target diagnostic codes or a prescription of an ADRD drug (See Tab. IV, Case Dx/Rx) |
| | **Control Cohort:** Patients lacking any target diagnostic code (Case Dx), or additionally any ADRD related prescription (Case Dx/Rx) |

CONSORT Diagram

1,033,782 patients 50+ $\geqslant$ 3yr record

→ 231,482 excluded insuff. history

729,018 used

387,700 female ↔ 341,318 male

TABLE II: Cohort Sizes

| | | Male | | | | Female | | |
|---|---|---|---|---|---|---|---|---|
| | **age range** | $n$ | $n_{\textbf{positive}}$ | $n_{\textbf{control}}$ | **age range** | $n$ | $n_{\textbf{positive}}$ | $n_{\textbf{control}}$ |
| **all patients** | 65-74 | 45943 | 1498 | 44445 | 65-74 | 47392 | 1546 | 45846 |
| | 75-84 | 16911 | 2870 | 14041 | 75-84 | 17303 | 3203 | 14100 |
| | 85+ | 7837 | 2383 | 5454 | 85+ | 10727 | 3753 | 6974 |
| | total | 341318 | 10397 | 330921 | total | 387700 | 12599 | 375101 |
| | **age range** | $n$ | $n_{\textbf{positive}}$ | $n_{\textbf{control}}$ | **age range** | $n$ | $n_{\textbf{positive}}$ | $n_{\textbf{control}}$ |
| **low-risk** [§] | 65-74 | 3841 | 108 | 3733 | 65-74 | 5330 | 101 | 5229 |
| | 75-84 | 963 | 170 | 793 | 75-84 | 1165 | 213 | 952 |
| | 85+ | 328 | 106 | 222 | 85+ | 411 | 170 | 241 |
| | total | 51499 | 797 | 50702 | total | 78152 | 1019 | 77133 |
| | **age range** | $n$ | $n_{\textbf{positive}}$ | $n_{\textbf{control}}$ | **age range** | $n$ | $n_{\textbf{positive}}$ | $n_{\textbf{control}}$ |
| **high-risk** [§] | 65-74 | 42102 | 1390 | 40712 | 65-74 | 42062 | 1445 | 40617 |
| | 75-84 | 15948 | 2700 | 13248 | 75-84 | 16138 | 2990 | 13148 |
| | 85+ | 7509 | 2277 | 5232 | 85+ | 10316 | 3583 | 6733 |
| | total | 289819 | 9600 | 280219 | total | 309548 | 11580 | 297968 |

[§]See Tab. I for cohort definitions, and SI-Table I for diagnostic codes defining the high-risk cohort.

inexpensive, fast and accessible solution, to accurately discover elevated risk of ADRD[21]. The multifactorial etiologies of ADRD imply that numerous risk factors are associated with these syndromes[21], and administrative claims and hospital databases, due to their large or even vast scope, offer sufficient statistical power to discover exquisitely-detailed algorithms to pinpoint potential cases. Notably, administrative claims and hospital databases may be especially amenable to exploration in heretofore unprecedented depth of associations of ADRD and comorbidities. Observational studies already have suggested that such associations encompass a large number and variety of disorders covering much of the human disease spectrum[22]; for example, non-neuropsychiatric chronic conditions such as diabetes, hypertension, hypercholesterolemia, obesity, sleep apnea, thyroid disorders, osteoporosis, and glaucoma have been linked to ADRD[22,23]. The validated or suspected associations of ADRD with both "intuitive" categories of comorbidity, $e.g.$, neurological, psychiatric, and cardiovascular disorders, and with non-intuitive categories, $e.g.$, metabolic, endocrine, opthalmologic disorders, and more recently infections[24,25], provide rationale for us to seek to leverage comorbid diagnoses to quantify ADRD risk.

Despite extensive documentation of co-morbidities, a reliable risk estimator — purely from ICD code-based co-morbidity patterns without any pre-selection of diagnostic codes already known to be a ADRD co-morbidity — is under-explored. The heterogeneity of brain aging processes and ADRD presentation[26], make such an endeavor challenging. Here we report the Zero Burden Co-morbid Risk Score (ZCoR) for ADRD, developed and validated using 729,018 unique patients drawn from across the United States, which reliably identifies patients up to 10 years before a contemporary documented clinical diagnosis.

ZCoR is a 701-feature digital signature distilled automatically from past diagnostic code sequences. We make no preselection of codes or ADRD-related risk factors, and require no new blood-work, laboratory tests, familial history or other patient-specific information that might preclude applicability at the point of care. Yet, we achieve an out-of-sample AUC exceeding 87% for either sex (predictions for one year earlier), and $\approx$ 80% (prediction

made a decade earlier). Importantly, our predictive performance matches the highest AUC achieved by MOCA (92.1%[19,27]) when making predictions just before a clinical diagnosis ($\geq$ 91%, see Table. IX). Our underlying algorithms are fundamentally novel, designed to learn from sparse, noisy categorical diagnostic sequences, with demonstrable non-trivial performance boost over standard tools used in recent studies.

Additionally, ZCoR for ADRD is sex-stratified, with separate signatures generated for males and females, in line with the growing appreciation of sex as a key contributor to the phenotypic heterogeneity of ADRD[28,29]. Predictability of eventual dementia in the presence of specific high-risk conditions such as type 2 diabetes has been studied before[30–32]. However, the complex pathobiology of ADRD implies that "low-risk" patients without any of the known risks, might still develop ADRD. Lacking the known flags, such a low-risk cohort is at a much higher risk of a missed or a delayed diagnosis. We show that ZCoR maintains high predictive performance in such patient groups.

The effectiveness of diverse cognitive assessment tools, often combined with pre-selected risk factors, has been recently surveyed[33], recording AUCs between 55%-89%, sometimes for diagnoses $10 - 20$ years into the future[34,35]. Nevertheless, substantial resource burden - often requiring detailed neurological, cognitive and psychiatric consults - limits applicability of such tools at the point-of-care, which were never conceived of for universal screening. Similarly, recent advances with machine learning[26] have often focused on classification of brain imaging data, which, while effective, and backed by well-understood mechanistic models, do not mitigate the barrier to universal adoption.

Thus, our key contribution in this study is to potentially alleviate obstacles to universal testing of the older population. The necessity of such universal screening tools has been well-recognized[36,37], with recent attempts at developing electronic health record (EHR)-based digital signatures to assess future ADRD risk. While two other digital signatures[36,37] have, to our knowledge, been reported since 2020, ZCoR demonstrates significantly better performance. More importantly, Boustani *et al.* used both structured and unstructured data including clinical notes processed for specific AD related keywords, and and Park *et al.* makes use of laboratory test results (*e.g.* blood hemoglobin) which might not be available for every patient at the point-of-care; in contrast, ZCoR exclusively uses data already present in patient records, which would typically vary from one patient to another, with no *a priori* fixed "demand" on any specific item of clinical, familial, demographic or lifestyle information. Thus, ZCoR can be applied almost universally, passively, and nearly instantaneously at the point-of-care.

EHR-based screening has also been explored to a limited extent for MCI[38], leveraging patterns extracted from clinical notes. When retrained to detect MCI, ZCoR-MCI is demonstrated to perform significantly better compared to reported results (both at both at the time of screening, and for predictions made years into the future), while, as before, using only ICD codes from past encounters.

## RESULTS

### Patient Selection

Our patient data comes from the IBM MarketScan® Commercial Claims and Encounters Database for the years 2003-2018[39] (previously Truven Health Analytics, and referred to as the "Truven dataset"). This US national database merges data contributed by over 150 insurance carriers and large self-insurance companies, and comprises over seven billion time-stamped diagnosis codes. The database tracks over 87 million patients for 1 to 15 years, reflecting a substantial cross-section of the US population. We select our cohort(s) in accordance with the inclusion/exclusion criteria described in Table I, ensuring that selected patients have at least three years of medical history recorded in the dataset. The geographical distribution of the patients in our selected cohort(s) is illustrated in Fig. 1a-b. Fig. 1c illustrates the age distribution at the time of ADRD diagnosis, which is consistent with the reported onset age characteristics for ADRD (mid-sixties[3]). Notably, the cumulative risk of onset (number of ADRD cases normalized by the total number of patients in the given age category) increases with age, as shown in Fig. 1c.

Predicting future ADRD diagnosis is modeled as a binary classification problem: we classify time-stamped sequences of diagnostic codes into positive and control categories, where the "positive" category refers to patients diagnosed with ADRD at 1 year from the point of screening, as identified by one or more ICD codes from Table III appearing in record (referred to as the Dx problem definition in Fig. 2 and Tables VII and VIII) or the prescription of AD-related medication[37,40] (Table IV, donepezil, galantamine, memantine or rivastigmine, referred to as the "Dx/Rx problem definition" in Table VII and VIII). The breakdown of diagnostic codes used for different sex and problem-definition combinations are shown in Table V.

We also consider screening up to $M$ years before the actual diagnosis, considering values of $M = 0, \cdots, 10$,

TABLE III: ADRD ICD diagnostic codes

| ICD code | description |
|---|---|
| 290.0 | Senile dementia uncomplicated |
| 290.1 | Presenile dementia |
| 290.11 | Presenile dementia w delirium |
| 290.12 | Presenile dementia w delusion |
| 290.13 | Presenile dementia w depression |
| 290.2 | Senile dementia w delusion |
| 290.21 | Senile dementia w depressive |
| 290.3 | Senile dementia w delirium |
| 290.8 | Senile psychosis NEC |
| 290.9 | Senile psychot condition NOS |
| 293.9 | Transient mental disease NOS |
| 294.2 | Demen NOS w/o behavior dstrb |
| 294.21 | Demen NOS w behavior distrb |
| 294.8 | Mental disorder NEC other disease |
| 294.9 | Mental disorder NOS other disease |
| 331.0 | Alzheimer's disease |
| F00 F00.0 F00.1 F00.2 F00.2 F00.9 | Dementia in Alzheimer's disease |
| F03.9 | Unspecified dementia without behavioral disturbance |
| F03.90 | Unspecified dementia without behavioral disturbance |
| F03.91 | Unspecified dementia with behavioral disturbance |
| F05 | Delirium due to known physiological condition |
| F06.8 | Other specified mental disorders due to known physiological condition |
| G30 | Alzheimer's disease with early onset |
| G30.0 | Alzheimer's disease with early onset |
| G30.1 | Alzheimer's disease with late onset |
| G30.8 | Other Alzheimer's disease |
| G30.9 | Alzheimer's disease unspecified |

TABLE IV: ADRD common prescriptions active ingredients

| Drugs |
|---|
| Donepezil Hydrochloride |
| Galantamine Hydrobromide |
| Memantine Hydrochloride |
| Rivastigmine |
| Rivastigmine Tartrate |

TABLE V: Number of diagnostic codes used

| target case | gender | Number of codes | Number of unique codes |
|---|---|---|---|
| Dx | Male | 6729970 | 18895 |
| Dx | Female | 9693456 | 19998 |
| Dx/Rx | Male | 6721267 | 18887 |
| Dx/Rx | Female | 9682339 | 19988 |

i.e., which relate to predicting ADRD immediately before a clinical diagnosis to up to a decade in the future. The control cohort comprises patients who never develop ADRD, i.e., do not have target codes, and are never prescribed related medication. Due to our requirement of minimum 3 years of medical history of record implies the absence of a diagnosis for at least $M + 2$ years in the future from the time of screening. We base our predictions on the past 2 years of diagnostic history. Overall we analyze $n = 729{,}018$ patients, with $22{,}996$ patients in the positive group and $706{,}022$ patients in the control group (See CONSORT diagram in Fig. 2c), considering approximately 42 million diagnostic codes, with altogether over 46K unique codes for both sexes.

We do not pre-select any diagnostic code based on its suspected comorbidity with ADRD. To investigate if our performance changes substantially for "high-risk" patients identified based on known co-morbidities including obesity, type II diabetes mellitus, hypertension, atherosclerosis, atrial fibrillation, dyslipidemia, depression, alcohol abuse, and pneumonia, we separately test our performance in high-risk and low-risk sub-cohorts. The high-risk sub-cohort comprises patients with one or more of the diagnoses enumerated in SI-Table I, which identify the top known co-morbidities[22,41]. The low-risk sub-cohort comprises patients who are not at high-risk as specified by the previous condition. Results in the low-risk sub-cohort is of particular significance; these patients are at a higher risk of missed or delayed diagnosis.

Determining the optimal set of target codes for making clinically useful predictions is challenging; too wide a definition makes predictions non-specific, while selecting too few erodes statistical power. The selection of target codes in Table III closely follows Park et al.[37] to enable a direct performance comparison. We also consider an expanded set of targets, including vascular dementia, frontotemporal dementia, vascular cognitive impairment, dementia with Lewy Bodies, and major neurocognitive disorder, with no significant performance variation (See Results).

The EHR codeset we use to ascertain ADRD-related disorders is intentionally broad, and not meant to diagnose a specific pathology, but predict risk of general dementia. This comports with our aim of developing a universal

TABLE VI: Feature Definitions (Total number of features used: 701)

| feature name | explanation | $n_{features}$ |
|---|---|---|
| **feature-phenotype** scores relative to phenotype score | Mean p-score of **feature-phenotype** codes within sequence divided by general p-score of **feature-phenotype** | 45 |
| **feature-phenotype** scores relative to whole score | Mean p-score of **feature-phenotype** codes within sequence divided by mean p-score of all codes in the record | 45 |
| aggregation score | aggregation of the p-scores in the record | 13 |
| high scores proportion | proportion of codes with very high p-scores among all codes in the record | 1 |
| low scores proportion | proportion of codes with very low p-scores among all codes in the record | 1 |
| dynamics of mean score | mean p-score of second half of the record divided by mean p-score of first half of the record | 1 |
| dynamics of geometric mean score | geometric mean p-score of second half of the record divided by mean p-score of first half of the record | 1 |
| dynamics of st.dev score | standard deviation of p-scores of second half of the record divided by standard deviation of p-scores of first half of the record | 1 |
| dynamics of score range | range of p-scores of second half of the record divided by range of p-scores of first half of the record | 1 |
| dynamics of score skew | skew of p-scores of seconf half of the record divided by skew of p-scores of first half of the record | 1 |
| aggregation relative to phn score | aggregation of all **feature-phenotype**'s mean scores divided by corresponding general p-score of **feature-phenotype** | 9 |
| aggregation relative to whole score | aggregation of all **feature-phenotype**'s mean scores divided by mean p-score of all codes in the record | 9 |
| **feature-phenotype** proportion | Ratio of number of weeks with the codes of a given phenotype to the total number of weeks in sequence | 45 |
| **feature-phenotype** prevalence | Ratio of number of weeks with the codes of a given phenotype to the number of weeks with any diagnosis code recorded | 45 |
| **feature-phenotype** first incident | Time interval from observation date to the first phenotype code, normalized by record length | 45 |
| **feature-phenotype** last incident | Time interval from observation date to the last phenotype code, normalized by record length | 45 |
| **feature-phenotype** mean position | Mean time position of phenotype codes in the record, normalized by record length | 45 |
| **feature-phenotype** streak | Length of the longest uninterrupted subsequence of weeks with the codes of a given phenotype recorded | 45 |
| **feature-phenotype** code prevalence | Ratio of number of codes of a given phenotype to the total number of codes in sequence | 45 |
| **feature-phenotype** code density | Ratio of number of codes of a given phenotype to the total number of weeks in sequence | 45 |
| Max/Mean/Std/Range intermission | Maximum/Mean/Standard Deviation/Range of the lengths of subsequences of consequent weeks with codes | 4 |
| Max/Mean/Std cluster | Maximum/Mean/Standard Deviation of the lengths of subsequences of consequent weeks without codes | 3 |
| Max/Std/Range prevalence | Maximum/Standard Deviation/Range of the phenotype prevalences | 3 |
| Max/Std/Range code prevalence | Maximum/Standard Deviation/Range of the Ratio of number of codes of a given phenotype to the total number of codes in sequence | 3 |
| Max/Std/Range code density | Maximum/Standard Deviation/Range of the Ratio of number of codes of a given phenotype to the total number of weeks in sequence | 3 |
| Density of DX Record | Proportion of weeks in a record observed where at least one DX code was recorded | 1 |
| **feature-phenotype** | Sequence Likelihood Defect for a given phenotype | 45 |
| **feature-phenotype** neg **log-likelihood** | negative log-likelihood score for a given phenotype | 45 |
| **feature-phenotype** pos **log-likelihood** | positive log-likelihood score for a given phenotype | 45 |
| **feature-phenotype log-likelihood** ratio | Ratio of positive to negative log-likelihood score for a given phenotype | 45 |
| Mean $\Delta$ ‡ | Mean negative Sequence Likelihood Defect | 1 |
| Geometric Mean $\Delta$ ‡ | Geometric Mean negative Sequence Likelihood Defect | 1 |
| Range $\Delta$ ‡ | Range of Sequence Likelihood Defect | 1 |
| Std. deviation $\Delta$ ‡ | Standard Deviation of Sequence Likelihood Defect | 1 |
| Mean neg **log-likelihood** | Mean negative log-likelihood score | 1 |
| Geometric Mean pos **log-likelihood** | Geometric Mean negative log-likelihood score | 1 |
| Range neg **log-likelihood** | Range of negative log-likelihood score | 1 |
| Std. deviation neg **log-likelihood** | Standard Deviation of negative log-likelihood score | 1 |
| Mean pos **log-likelihood** | Mean positive log-likelihood score | 1 |
| Geometric Mean pos **log-likelihood** | Geometric Mean positive log-likelihood score | 1 |
| Range pos **log-likelihood** | Range of positive log-likelihood score | 1 |
| Std. deviation pos **log-likelihood** | Standard Deviation of positive log-likelihood score | 1 |
| Mean **log-likelihood** ratio | Mean log-likelihood score ratio | 1 |
| Geometric Mean **log-likelihood** ratio | Geometric Mean log-likelihood score ratio | 1 |
| Range **log-likelihood** ratio | Range of log-likelihood score ratio | 1 |
| Std. deviation **log-likelihood** ratio | Standard Deviation of log-likelihood score ratio | 1 |

*feature: ICD disease categories, or sets of diagnostic codes tracked
†$\Delta$: Sequence Likelihood Defect (See Methods)
‡ neg log-likelihood: log-likelihood of observed sequence generated by model inferred from control (See Methods)
# pos log-likelihood: log-likelihood of observed sequence generated by model inferred from positive (See Methods)

**a.** Receiver Operator Characteristic curves

**b.** Precision-Recall curves

**c.** Feature importances for broad categories of co-morbidities

Fig. 2: **Predictive performance of ZCoR for ADRD diagnosis 1 year in the future.** Panels a and b show the out-of-sample ROC and precision-recall curves for diagnosis 1 year from the point of screening. We achieve AUCs $> 88\%$ for male and $> 86\%$ for females in the age group 50+, for the diagnostic criteria based on ICD codes (See description of diagnostic criteria considered in Table I), with sensitivities at 58% (females) and 54% (females) at 95% specificity. See Tables VII and VIII for performance within 65+ cohort, and within the low-risk and high-risk cohorts in each age strata. Panel c shows the top 20 comorbidity categories sorted in the order of inferred importance in estimating risk, where categories for mental and cognitive disorders have been removed to highlight the role of other physiological co-morbidities. Importantly, the comorbidities modulate risk differentially by sex, although the patterns are broadly similar, *e.g.*, metabolic, cardiovascular, opthalmological, ischemic categories appear in both males and females, with slightly altered ranking. Infections and immunologic disorders appear with high importance.

screening tool, as opposed to a diagnostic intrument, that triggers more detailed neurological assessment.

In addition to ADRD, we also investigate the ability of our basic approach to predict MCI (both at the time of screening, and as a prediction of a future diagnosis), identified by the appearance of ICD codes 331.83 (ICD9) and G31.84 (ICD10). We evaluate the performance of ZCoR-MCI with a retrained pipeline, which, in out-of-sample validation, shows significant improvement over reported literature.

*Feature Importance & Comorbidity Spectra*

The aggregate importance of the ZCoR features (See Fig. 2c), estimated as the mean change in the raw risk via random perturbations in the feature values, illustrates that metabolic and cardiovascular disorders are the most important diagnostic category modulating risk.

Additionally, we compute the statistically significant log-odds ratio of specific ICD codes occurring in the true positive vs the true negative patient sets. We call these the "comorbidity spectra" (See Figs. 3 and 4). These spectra are based on individual codes, as opposed to the aggregated feature importances shown in Fig. 2c. Clearly, every disorder listed in the co-morbid spectra does not all appear in a single patient; the codes with high log-odds ratio are significantly more likely in the positive cohort. The comorbidity spectra, so named because of disease category-specific color coding, offers unique insight into the predictive co-morbidity burden of ADRD.

## Outcomes

In this study we demonstrate the following key results: 1) high out-of-sample predictive performance for identifying a ADRD diagnosis 1 year into future via leveraging subtle comorbidity patterns recorded in the medical history of individual patients (Tables VII-VIII), 2) high predictive performance for diagnosis up to 10 years into the future with sufficiently slow loss of predictive performance to remain clinically useful (Table IX), significantly outperforming recent results (Tables X-XI), 3) effective performance for both low-risk and high-risk cohorts (Tables VII-VIII, see relevant rows). Here, the high-risk cohort comprises patients with commonly surveilled for ADRD co-morbidities. Additionally, 4) maintain high performance for expanded target definitions which include vascular dementia, frontotemporal dementia, vascular cognitive impairment, dementia with Lewy Bodies, and major neurocognitive disorders (Table XII). And finally, 5) high predictive performance to screen for MCI for current and future diagnosis (Table XIII).

Pertianing to our main prediction results for 1-3, Fig. 2a-b illustrate the ROC and the precision-recall curves respectively (for screening one year before current diagnosis), shown separately for males and females. As noted in the panel legends, our out-of-sample predictive performance is $> 88\%$ AUC for females (age 50+) and $> 86\%$ for males (age 50+), with $> 50\%$ sensitivity at 95% specificity (53% for males and 57% for females). At 99% specificity, we obtain a PPV of 42% for females (50+) and $40 - 41\%$ for males (50+) respectively. At these values we obtain an accuracy of $\approx 96 - 97\%$ (Table VII), which indicates the overall fraction of correct predictions. The PPV achieved by ZCoR at maximum accuracy is $54 - 55\%$ for females (50+) and $51 - 53\%$ for males (50+), with a corresponding NPV of 97%. The corresponding results for age 65+ are tabulated in Table VIII.

Thus, to summarize: our predictive pipeline detects about 53-57 out of every 100 patients who get a diagnosis in 1 year, if we operate at 95% specificity. If we wish to operate at the higher specificity of 99%, then out of 100 positive flags, we have about 41-42 true positives. The accuracy metric indicates that we correctly identify the risk status (positive or control) for 96-97 out of 100 patients, irrespective of sex, highlighting the potentially high clinical significance of ZCoR, as a universal screening tool to identify patients for diagnostic workup and/or intensified surveillance.

From the inferred relative importance of the co-morbidity categories (See Fig. 2d-e), we conclude, that metabolic and ischemic diseases, cardiovascular abnormalities, sleep disorders, nervous system disorders, and diseases of the eye are important modulators of risk. Infections also feature in the top 20 co-morbidities shown in these panels. Importantly while there are sex differences, the overall pattern of the relative importance ranking remains substantially sex-invariant. With some exceptions, many of these patterns are not particularly surprising; the contribution of this study is to bring them together systematically to realize an accurate risk estimate via the ZCoR score.

As expected, our predictive performance degrades as we predict earlier (See Table IX, and inset). Importantly, the degradation is slow enough that we can use ZCoR with acceptable reliability up to 10 years into the future, and significantly outperform reported results.

Understanding the seat of this predictive power is important. The feature importances discussed earlier (Fig. 2c) identify the relative impact of broad disease categories. Importantly, to evaluate the feature importance of a specific diagnostic category, we sum the importance of all features based on that category, not just the presence or absence of individual diagnoses. The latter aspect, $i.e.$, the risk burden from the presence of specific codes, is investigated via the co-morbidity spectra for out-of-sample patients, shown separately in Figs. 3 and 4 for males and females, and the two target definitions (Dx and Dx/Rx). We find that the important co-morbidities are diverse, vary with the sex of the patients, but are clearly dominated by mental disorders, circulatory disorders, injuries, and a range of disorders categorized broadly as "ill-defined symptoms" in the ICD framework. Again,

TABLE VII: Detailed ZCoR performance for patients aged 50+, predictions made 1 year before diagnosis

| sex | definition | cohort | sens. | PPV | acc | PPV[†] | NPV[†] | spec. | auc |
|---|---|---|---|---|---|---|---|---|---|
| Female | Dx/Rx | all patients | 0.57 | 0.29 | 0.94 | 0.55 | 0.97 | 95% | $0.884 \pm 0.008$ |
| Female | Dx | all patients | 0.58 | 0.29 | 0.94 | 0.54 | 0.97 | 95% | $0.885 \pm 0.004$ |
| Female | Dx/Rx | all patients | 0.22 | 0.42 | 0.96 | 0.55 | 0.97 | 99% | $0.884 \pm 0.008$ |
| Female | Dx | all patients | 0.21 | 0.42 | 0.96 | 0.54 | 0.97 | 99% | $0.885 \pm 0.004$ |
| Male | Dx/Rx | all patients | 0.53 | 0.25 | 0.94 | 0.51 | 0.97 | 95% | $0.866 \pm 0.006$ |
| Male | Dx | all patients | 0.54 | 0.27 | 0.94 | 0.53 | 0.97 | 95% | $0.870 \pm 0.011$ |
| Male | Dx/Rx | all patients | 0.21 | 0.41 | 0.97 | 0.51 | 0.97 | 99% | $0.866 \pm 0.006$ |
| Male | Dx | all patients | 0.21 | 0.40 | 0.97 | 0.53 | 0.97 | 99% | $0.870 \pm 0.011$ |
| Female | Dx/Rx | high-risk | 0.55 | 0.28 | 0.94 | 0.54 | 0.97 | 95% | $0.883 \pm 0.008$ |
| Female | Dx | high-risk | 0.55 | 0.28 | 0.94 | 0.50 | 0.97 | 95% | $0.883 \pm 0.005$ |
| Female | Dx/Rx | high-risk | 0.20 | 0.41 | 0.96 | 0.54 | 0.97 | 99% | $0.883 \pm 0.008$ |
| Female | Dx | high-risk | 0.19 | 0.40 | 0.96 | 0.50 | 0.97 | 99% | $0.883 \pm 0.005$ |
| Male | Dx/Rx | high-risk | 0.52 | 0.25 | 0.94 | 0.58 | 0.97 | 95% | $0.867 \pm 0.008$ |
| Male | Dx | high-risk | 0.53 | 0.25 | 0.94 | 0.50 | 0.97 | 95% | $0.871 \pm 0.011$ |
| Male | Dx/Rx | high-risk | 0.21 | 0.40 | 0.97 | 0.58 | 0.97 | 99% | $0.867 \pm 0.008$ |
| Male | Dx | high-risk | 0.20 | 0.39 | 0.97 | 0.50 | 0.97 | 99% | $0.871 \pm 0.011$ |
| Female | Dx/Rx | low-risk | 0.59 | 0.28 | 0.94 | 0.64 | 0.98 | 95% | $0.830 \pm 0.039$ |
| Female | Dx | low-risk | 0.54 | 0.28 | 0.94 | 0.62 | 0.98 | 95% | $0.833 \pm 0.039$ |
| Female | Dx/Rx | low-risk | 0.41 | 0.58 | 0.97 | 0.64 | 0.98 | 99% | $0.830 \pm 0.039$ |
| Female | Dx | low-risk | 0.37 | 0.56 | 0.97 | 0.62 | 0.98 | 99% | $0.833 \pm 0.039$ |
| Male | Dx/Rx | low-risk | 0.54 | 0.26 | 0.94 | 0.56 | 0.97 | 95% | $0.816 \pm 0.039$ |
| Male | Dx | low-risk | 0.59 | 0.28 | 0.94 | 0.68 | 0.97 | 95% | $0.826 \pm 0.029$ |
| Male | Dx/Rx | low-risk | 0.28 | 0.48 | 0.97 | 0.56 | 0.97 | 99% | $0.816 \pm 0.039$ |
| Male | Dx | low-risk | 0.31 | 0.51 | 0.97 | 0.68 | 0.97 | 99% | $0.826 \pm 0.029$ |

*Calculated at 95% specificity
[†]Maximum PPV at observed prevalence, and NPV at maximum PPV

TABLE VIII: Detailed ZCoR performance for patients aged 65+, predictions made 1 year before diagnosis

| sex | definition | cohort | sens. | PPV | acc | PPV[†] | NPV[†] | spec. | auc |
|---|---|---|---|---|---|---|---|---|---|
| Female | Dx/Rx | all patients | 0.29 | 0.43 | 0.87 | 0.56 | 0.89 | 95% | $0.836 \pm 0.007$ |
| Female | Dx | all patients | 0.28 | 0.42 | 0.87 | 0.54 | 0.89 | 95% | $0.832 \pm 0.010$ |
| Female | Dx/Rx | all patients | 0.09 | 0.54 | 0.89 | 0.56 | 0.89 | 99% | $0.836 \pm 0.007$ |
| Female | Dx | all patients | 0.07 | 0.49 | 0.88 | 0.54 | 0.89 | 99% | $0.832 \pm 0.010$ |
| Male | Dx/Rx | all patients | 0.31 | 0.41 | 0.89 | 0.54 | 0.91 | 95% | $0.827 \pm 0.011$ |
| Male | Dx | all patients | 0.31 | 0.40 | 0.89 | 0.54 | 0.91 | 95% | $0.829 \pm 0.013$ |
| Male | Dx/Rx | all patients | 0.10 | 0.53 | 0.90 | 0.54 | 0.91 | 99% | $0.827 \pm 0.011$ |
| Male | Dx | all patients | 0.09 | 0.51 | 0.90 | 0.54 | 0.91 | 99% | $0.829 \pm 0.013$ |
| Female | Dx/Rx | high-risk | 0.26 | 0.43 | 0.87 | 0.57 | 0.89 | 95% | $0.831 \pm 0.007$ |
| Female | Dx | high-risk | 0.27 | 0.42 | 0.87 | 0.57 | 0.89 | 95% | $0.827 \pm 0.009$ |
| Female | Dx/Rx | high-risk | 0.09 | 0.54 | 0.88 | 0.57 | 0.89 | 99% | $0.831 \pm 0.007$ |
| Female | Dx | high-risk | 0.07 | 0.48 | 0.88 | 0.57 | 0.89 | 99% | $0.827 \pm 0.009$ |
| Male | Dx/Rx | high-risk | 0.31 | 0.41 | 0.89 | 0.53 | 0.91 | 95% | $0.824 \pm 0.013$ |
| Male | Dx | high-risk | 0.30 | 0.40 | 0.89 | 0.52 | 0.91 | 95% | $0.827 \pm 0.014$ |
| Male | Dx/Rx | high-risk | 0.10 | 0.53 | 0.90 | 0.53 | 0.91 | 99% | $0.824 \pm 0.013$ |
| Male | Dx | high-risk | 0.09 | 0.51 | 0.90 | 0.52 | 0.91 | 99% | $0.827 \pm 0.014$ |
| Female | Dx/Rx | low-risk | 0.52 | 0.58 | 0.90 | 0.64 | 0.92 | 95% | $0.888 \pm 0.025$ |
| Female | Dx | low-risk | 0.45 | 0.54 | 0.89 | 0.55 | 0.93 | 95% | $0.869 \pm 0.033$ |
| Female | Dx/Rx | low-risk | 0.17 | 0.71 | 0.89 | 0.64 | 0.92 | 99% | $0.888 \pm 0.025$ |
| Female | Dx | low-risk | 0.10 | 0.90 | 0.89 | 0.55 | 0.93 | 99% | $0.869 \pm 0.033$ |
| Male | Dx/Rx | low-risk | 0.41 | 0.48 | 0.90 | 0.61 | 0.92 | 95% | $0.866 \pm 0.050$ |
| Male | Dx | low-risk | 0.37 | 0.46 | 0.89 | 0.57 | 0.92 | 95% | $0.858 \pm 0.037$ |
| Male | Dx/Rx | low-risk | 0.15 | 0.62 | 0.91 | 0.61 | 0.92 | 99% | $0.866 \pm 0.050$ |
| Male | Dx | low-risk | 0.14 | 0.64 | 0.91 | 0.57 | 0.92 | 99% | $0.858 \pm 0.037$ |

*Calculated at 95% specificity
[†]Maximum PPV at observed prevalence, and NPV at maximum PPV

while many of these patterns are known at the population level, design of the personalized ZCoR score is not immediately obvious.

We include predictive performance in conventional high-risk (defined in SI-Table I) and low-risk cohorts in Table VII and VIII, showing that our performance in the high-risk sub-cohort is comparable with that in the full cohort. The AUCs in the low-risk cohort are somewhat lower ($> 81\%$ for males 50+ and $> 83\%$ for females 50+ respectively), albeit high enough to be clinically effective: we have a maximum PPV of $62 - 68\%$, and a sensitivity of $54 - 59\%$ at specificity of 95% for 50+ patients who get diagnosed 1 year in the future (See

TABLE IX: Long-range ZCoR AUC estimates (95% confidence bounds) for target set listed in Table III

| years to diagnosis | AUC Female | AUC Male | AUC Female Dx/Rx* | AUC Male Dx/Rx |
|---|---|---|---|---|
| 0 | $0.912 \pm 0.015$ | $0.913 \pm 0.015$ | $0.909 \pm 0.015$ | $0.918 \pm 0.015$ |
| 1 | $0.885 \pm 0.015$ | $0.871 \pm 0.015$ | $0.875 \pm 0.015$ | $0.867 \pm 0.015$ |
| 2 | $0.872 \pm 0.017$ | $0.858 \pm 0.017$ | $0.867 \pm 0.017$ | $0.858 \pm 0.017$ |
| 3 | $0.858 \pm 0.019$ | $0.850 \pm 0.018$ | $0.852 \pm 0.019$ | $0.855 \pm 0.019$ |
| 4 | $0.863 \pm 0.021$ | $0.842 \pm 0.020$ | $0.843 \pm 0.021$ | $0.860 \pm 0.021$ |
| 5 | $0.850 \pm 0.023$ | $0.841 \pm 0.022$ | $0.840 \pm 0.024$ | $0.846 \pm 0.023$ |
| 6 | $0.840 \pm 0.026$ | $0.831 \pm 0.025$ | $0.832 \pm 0.027$ | $0.835 \pm 0.025$ |
| 7 | $0.830 \pm 0.031$ | $0.810 \pm 0.028$ | $0.839 \pm 0.031$ | $0.835 \pm 0.028$ |
| 8 | $0.815 \pm 0.036$ | $0.815 \pm 0.033$ | $0.828 \pm 0.037$ | $0.823 \pm 0.034$ |
| 9 | $0.799 \pm 0.046$ | $0.809 \pm 0.041$ | $0.811 \pm 0.049$ | $0.807 \pm 0.042$ |
| 10 | $0.841 \pm 0.059$ | $0.784 \pm 0.049$ | $0.834 \pm 0.062$ | $0.810 \pm 0.058$ |

* Dx/Rx refers to diagnosis inferred from either codes or AD-related prescriptions

**INSET. ZCoR AUC over time**



TABLE X: Comparison of AUC achieved in out-of-sample data between ZCoR and Park et al.[37]

| Year to diagnosis | Park et al. (Dx/Rx) | Park et al. (Dx) | ZCoR (Dx) | ZCoR (Dx/Rx) | Δ(Dx)%[‡] | Δ(Dx/Rx)%[†] |
|---|---|---|---|---|---|---|
| 0 | 0.90 | 0.85 | 0.91 | 0.92 | 7.4726 | 2.3315 |
| 1 | 0.78 | 0.76 | 0.89 | 0.88 | 16.621 | 12.933 |
| 2 | 0.73 | 0.69 | 0.87 | 0.87 | 25.851 | 18.793 |
| 3 | 0.68 | 0.64 | 0.86 | 0.86 | 33.373 | 26.404 |
| 4 | 0.72 | 0.68 | 0.86 | 0.86 | 26.467 | 18.639 |

TABLE XI: Comparison of AUC achieved in out-of-sample data between ZCoR and Boustani et al.[36]

| year to diagnosis | Boustani et al. (Dx) | ZCoR (Dx) | Δ(Dx)%[‡] |
|---|---|---|---|
| 1-10 | 0.80 | 0.85 | 6.48 |
| 3-10 | 0.75 | 0.84 | 11.9 |
| 5-10 | 0.70 | 0.83 | 17.8 |

[‡] Percentage outperformance of ZCoR with the Dx target definition
[†] Percentage outperformance of ZCoR with the Dx/Rx target definition

Tables VII and VIII).

Additionally, Table XII shows that an expanded target definition (described in Methods) yields no significant change in predictive performance. Finally, our results on MCI prediction are shown in Table XIII (carried out with a retrained pipeline targeting MCI), which illustrates an average AUC between 88-90% at the point of screening, degrading to under $\approx 80\%$ for predictions made 3 years into the future. Notably, the performance at the point of screening is at par with MOCA ($0.9 - 0.91 \pm 0.015$ vs $.921$[19]), while being significantly superior to a recently published EHR-based approach using standard machine learning algorithms[42].

## DISCUSSION

We report the development and validation of the ZCoR automated universal screening tool for ADRD, leveraging previously-uncharted co-morbidity patterns discovered from individual longitudinal diagnostic history. Across sexes, ZCoR accurately preempts ADRD cases up to 10 years before a clinical diagnosis is first documented. The broad co-morbidity categories that we infer to be important (Fig. 2c) include metabolic, cardiovascular, ischemic, ophthalmological, and sleep disorders. Diseases of the nervous system, unrelated to ADRD, infections, and immunologic disorders also appear in the list of top risk-modulating co-morbidities. Importantly, the co-morbidities modulate risk differentially by sex (Fig. 2c), although the patterns are broadly similar, with slightly altered ranking.

TABLE XII: Long-range ZCoR$^\star$ AUC estimates (95% confidence bounds) for expanded target set (adding vascular dementia, frontotemporal dementia, vascular cognitive impairment, dementia with Lewy Bodies, and major neurocognitive disorder to Table III).

| years to diagnosis | AUC Female | AUC Male |
|---|---|---|
| 0 | $0.907 \pm 0.155$ | $0.913 \pm 0.155$ |
| 1 | $0.884 \pm 0.155$ | $0.871 \pm 0.155$ |
| 2 | $0.876 \pm 0.173$ | $0.858 \pm 0.170$ |
| 3 | $0.871 \pm 0.191$ | $0.853 \pm 0.187$ |
| 4 | $0.857 \pm 0.210$ | $0.847 \pm 0.206$ |
| 5 | $0.846 \pm 0.236$ | $0.836 \pm 0.226$ |
| 6 | $0.833 \pm 0.262$ | $0.822 \pm 0.253$ |
| 7 | $0.821 \pm 0.299$ | $0.829 \pm 0.294$ |
| 8 | $0.807 \pm 0.365$ | $0.804 \pm 0.358$ |
| 9 | $0.787 \pm 0.473$ | $0.769 \pm 0.456$ |
| 10 | $0.793 \pm 0.606$ | $0.769 \pm 0.640$ |

**INSET.** ZCoR$^\star$ AUC over time



TABLE XIII: Long-range ZCoR-MCI AUC estimates (95% confidence bounds) with MCI identified via ICD codes 331.83 (ICD9) and G31.84 (ICD10).

| years to diagnosis | AUC Female | AUC Male |
|---|---|---|
| 0 | $0.882 \pm 0.038$ | $0.901 \pm 0.040$ |
| 1 | $0.845 \pm 0.038$ | $0.850 \pm 0.040$ |
| 2 | $0.846 \pm 0.042$ | $0.828 \pm 0.043$ |
| 3 | $0.818 \pm 0.046$ | $0.789 \pm 0.046$ |
| 4 | $0.780 \pm 0.050$ | $0.821 \pm 0.050$ |
| 5 | $0.783 \pm 0.054$ | $0.805 \pm 0.055$ |
| 6 | $0.758 \pm 0.060$ | $0.790 \pm 0.061$ |
| 7 | $0.764 \pm 0.066$ | $0.800 \pm 0.069$ |
| 8 | $0.747 \pm 0.079$ | $0.792 \pm 0.080$ |
| 9 | $0.694 \pm 0.102$ | $0.765 \pm 0.100$ |
| 10 | $0.611 \pm 0.168$ | $0.764 \pm 0.132$ |

**INSET.** ZCoR-MCI AUC over time



Focusing on the presence/absence of individual diagnostic codes modulating ADRD risk in the co-morbidity spectra (Figs. 3 and 4), we find circulatory disorders are generally over-represented, along with injuries, and conditions related to age-related cognitive decline. Many of these patterns are unsurprising: AD is a amnestic syndrome, injuries might indicate neuropathies from known AD co-morbidities such as diabetes or stroke, and cerebrovascular diseases might signal vascular dementia. Other prominent codes such as ataxia and psychiatric signs were recently associated with specific biomarkers implicated in autosomal dominant early-onset Alzheimer's disease[43,44]. Appearance of other codes are more surprising, $e.g.$ dysphagia or swallowing impairment is usually noted in the late stages of AD. However, recent studies have documented changes in cortical control of swallowing beginning before dysphagia becomes apparent in dementia patients[45,46]. Thus, the important illness categories that we find to be associated with ADRD in either sex align with suspected or documented links in statistical[29,36,37] and observational studies[22,23], lending credence to ZCoR rationale and accuracy. Also lending such credence is the score's incorporation, via sex-stratification, of differences between males and females in ADRD risk factors, natural history, and symptoms[28,29,47–51].

We find that with increasing patient age, it becomes more difficult to distinguish age related cognitive decline from ADRD (SI-Fig. 1), suggesting that ADRD comorbidities have confounding overlaps with conditions that arise more frequently as patients get older.

To our knowledge, ZCoR is one of three digital signatures for ADRD reported since 2020, joining those of Boustani $et\,al.$, developed utilizing data from the Indiana Network for Patient Care[36], and of Park $et\,al.$, developed utilizing data from the Korean National Health Insurance Service[37]. Although the respective reported prognostic time-frames are not fully comparable, our digital signature appeared to achieve the best performance of the three (Table IX inset, and Tables X and XI). Notably, the AUC of ZCoR for ADRD at 10 years before documented diagnosis surpassed the AUCs of the Boustani $et\,al.$[36] signature for the 1-10 year, 3-10 year, or 5-10 year before diagnosis time-frames by 6.5%, 11.9%, and 17.8% respectively, while leveraging diagnostic histories of 1,400% more patients ($\approx$ 50K vs $\approx$ 700K for ZCoR). Also for each prediction time-point made 0 through 4 years before documented diagnosis, the AUCs of ZCoR exceeded those of the Park $et\,al.$[37] signature by 2-7% (0

## ICD Class

- Infections
- Endocrine & Immun. Dis.
- Skin & subcut. tiss.
- Eye & adnexa
- Ear & mastoid
- Circulatory Dis.
- Blood and bld. form. organs
- Musculosk. & Conn. Tiss.
- Digestive Dis.
- Respiratory Dis.
- Genitourinary
- Neoplasms
- Mental Dis.
- Nervous Dis.
- Congenital Anomaly
- Cond. orig. in Perinatal Per.
- Injury & Poisoning
- External morbidity
- Ill-defined Cond. & Symp.
- Health service Contact

### a. Male Alzheimer's Dis.

| Code | Description |
|---|---|
| 518.5 | Ac resp flr fol trma/srg |
| I66.9 | Occlusion stenosis |
| S88.1 | Traumatic amput knee ankle lower leg |
| N32.9 | Bladder dis non-sp |
| R26.0 | Ataxic gait |
| C34.1 | Malig neopl up lobe non-sp bronchus lung |
| N32.3 | Diverticulum bladder |
| Z91.8 | Hisry falling |
| 780.6 | Fever nos |
| D51.0 | Vitamin b12 deficiency anemia intrinsic fact deficiency |
| 596.8 | Inf cyssmy |
| I44.3 | Non-sp atrioventricular block |
| I63.3 | Cerebral infarction thrombosis |
| 173.0 | Malig neopl skin lip nos |
| 173.6 | Mal neo skin up limb nos |
| 599.7 | Hematuria nos |
| N31.9 | Neuromuscular dysfunction bladder non-sp |
| I44.2 | Atrioventricular block complete |
| I95.1 | Orthostatic hypotension |
| I44.1 | Atrioventricular block second |
| I63.4 | Cerebral infarction embolism |
| S41.0 | Wound rt shoulder |
| I62.0 | Nontraumatic subdural hemorrhage non-sp |
| I69.9 | Nonsp sequelae cerebrovascular dis |
| I67.8 | Ac cerebrovascular insufficiency |
| I63.5 | Cerebral infarction non-sp occlusion stenosis |
| R47.0 | Aphasia |
| I67.1 | Cerebral aneurysm nonruptured |
| F06.3 | Mood dis known physcondition non-sp |
| I95.2 | Hypotension drugs |
| C67.8 | Malig neopl overlappings bladder |
| D41.4 | Neoplasm uncertain behavi bladder |
| D30.3 | 30-39pc bdy brn/30-39pc 3d |
| R40.4 | Transient alteration awareness |
| R27.0 | Ataxia non-sp |
| 173.2 | Malig neo skin ear nos |
| I61.9 | Nontraum intracerebral hem non-sp |
| R41.8 | Age-related cognitive decline |
| 787.2 | Dysphagia, oral phase |
| F29 | Non-sp psychosis |
| R41.3 | Amnesia |
| F05 | Delirium known physcondition |

log odds ratio of normalized prevalence

### b. Female Alzheimer's Dis.

| Code | Description |
|---|---|
| R26.0 | Ataxic gait |
| Z99.8 | Dependence on supplemental oxygen |
| S59.1 | Non-sp physeal fracture up end radius rt arm |
| S22.0 | Wedge compression fracture non-sp thoracic vertebra |
| S42.0 | Fracture non-sp part rt clavicle |
| Z45.0 | Checking testing cardiac pacemaker pulse generat [battery] |
| S30.0 | Contusion lower back pelvis |
| H35.3 | Non-sp macular degeneration |
| 787.2 | Dysphagia, oral phase |
| I50.1 | Left ventricular failure non-sp |
| C19 | Malig neopl recsigmoid junction |
| D04.3 | Carcinoma in situ skin non-sp part face |
| S70.0 | Contusion non-sp hip |
| S09.9 | Non-sp inj head |
| R40.0 | Somnolence |
| I50.9 | Heart failure non-sp |
| I74.3 | Embolism thrombosis arteries lower extremities |
| I95.1 | Orthostatic hypotension |
| I69.9 | Nonsp sequelae cerebrovascular dis |
| S42.3 | Non-sp fracture shaft humerus rt arm |
| 814.0 | Fx navicular wrist-clos |
| S39.8 | Sp injuries abdomen |
| 173.9 | Malig neo skin nos |
| R63.0 | Anorexia |
| J82 | Chronic eosinophilic pneumonia |
| S02.2 | Fracture nasal bones |
| I48.9 | Non-sp atrial fibrillation |
| L97.9 | Non-pressure chronic ulcer non-sp lower leg |
| K26.9 | Duodenal ulcer non-sp acute chronic wo hemorrhage perf |
| J80 | Ac respiray distress syndrome |
| I80.3 | Phlebitis thrombophlebitis lower extremities non-sp |
| I67.2 | Cerebral atherosclerosis |
| 829.0 | Fracture nos-closed |
| S01.0 | Wound scalp |
| I45.9 | Conduction dis non-sp |
| I49.5 | Sick sinus syndrome |
| S01.8 | Wound other part head |
| N36.2 | Urethral caruncle |
| I69.8 | Non-sp sequelae cerebrovascular dis |
| I65.8 | Occlusion stenosis other precerebral arteries |
| Z95.0 | Presence cardiac pacemaker |
| S32.0 | Wedge compression fracture non-sp lumbar vertebra |
| I44.2 | Atrioventricular block complete |
| 173.3 | Mal neo skn face nec/nos |
| I67.9 | Cerebrovascular disease non-sp |
| J69.0 | Pneumonitis inhalation food vomit |
| J81.1 | Chronic pulmonary edema |
| 813.8 | Fx radius nos-closed |
| 453.8 | Ac embl suprfcl up ext |
| S79.8 | Sp injuries rt hip |
| M80.0 | Age-related osteoporosis |
| R27.0 | Ataxia non-sp |
| R41.8 | Age-related cognitive decline |
| I50.4 | Non-sp congestive heart failure |
| H02.1 | Non-sp ectropion rt up eyelid |
| F29 | Non-sp psychosis |
| R40.4 | Transient alteration awareness |
| I63.5 | Cerebral infarction non-sp occlusion stenosis |
| I67.8 | Ac cerebrovascular insufficiency |
| 820.0 | Fx up femur epiphy-clos |
| R41.3 | Amnesia |
| 820.8 | Fx neck femur nos-cl |

ICD10 codes

log odds ratio of normalized prevalence

**Fig. 3: Co-morbidity Spectrum for the Dx/Rx case.** Disorders that increase the odds of the patient being a "true positive" vs a "true negative", where diagnosis is determined using either ICD codes (See Table III) or ADRD-related medications (See Table IV) in history. Such disorders (ranked according to the log-odds ratio) are more likely to be found in patients who are in the positive cohort. Comapring **panel a** with **panel b**, we note that these odds change from males to females, but as expected the patterns are broadly similar, with over-representation of circulatory disorders.

year), 12.9-16.6% (1 year), 18.8-25.8% (2 years), 26.4-33.3% (3 years) and 18.6-26.5% (4 years), while using 1,750% more patients ($\approx$ 40K vs $\approx$ 700K for ZCoR).

A number of methodological differences contribute to these respective performances. First, ZCoR uses sophisticated pattern discovery on patient history, and is not limited by known risk factors and co-morbidities, allowing for high performance on low-risk and high-risk cohorts alike. More importantly, our new stochastic inference

**Fig. 4: Co-morbidity Spectrum for the Dx case.** Disorders that increase the odds of the patient being a "true positive" vs a "true negative"', where diagnosis is determined using ICD codes (See Table III). Such disorders (ranked according to the log-odds ratio) are more likely to be found in patients who are in the positive cohort. Comapring **panel a** with **panel b**, we note that these odds change from males to females, but as expected the patterns are broadly similar, with over-representation of circulatory disorders.

algorithms are designed to leverage longitudinal patterns, and are not limited to using indicator variables, *i.e..*, simply the presence or absence of specific historical codes. Thus we are able to substantially leverage the emergent dependencies and temporal ordering of patterns emergent across the human disease spectrum.

Additionally, ZCoR for ADRD is stratified by sex. Sex-stratification of AD risk has recently found support in the literature[29]. Finally, our algorithm is derived using a cohort roughly 10-18-fold larger than those of Boustani

*et al.* or Park *et al.* (729,018 versus 40,736 and 71,466 respectively), allowing our algorithms to capitalize on significantly larger quantities of data.

Beyond predictive performance, ZCoR addresses the barrier to universal testing. With no specific data demands (we use what we have on the individual patient file), and designed to operate on existing electronic healthcare architectures, the digital signature operates non-invasively, inexpensively, and nearly instantaneously, and is potentially very widely, if not universally accessible at least in developed countries using EHR. Unlike that of Boustani *et al.*, (but like that of Park *et al.*) who use expert opinion-generated variables in the first phase of their digital signature development, our algorithm is completely data-driven. Also, unlike Boustani *et al.* (but like Park *et al.*), we considered only structured data, *i.e.*, ICD codes, and not clinical notes. While clinical notes might reveal substantially more information, such insights most relevant to ADRD might not be available before a neurology consult. And unlike Park *et al.*, we do not use laboratory tests such as hemoglobin levels, which might not be available for every patient in primary care.

We envision three main potential applications of ZCoR. First, the score can serve in primary care or specialist settings (*e.g.*, neurology, gerontology) as a screening tool for future incident overt cases, with the potential diagnostic, therapeutic, psychosocial, caregiver-related, and research benefits noted in the introduction to this paper. ZCoR could, for example, be routinely deployed, alone or along with a brief, validated neuropsychological instrument, as recommended by the American Academy of Neurology[52], in the cognitive screening mandated since 2011 as part of the Medicare annual wellness visit[3]. Alternatively, especially given the variable clinical natural history of such patients[53], ZCoR could be employed in individuals with subjective memory decline but largely-intact cognition and function, or in those with incipient MCI, *e.g.*, worsening but still personally-appropriate serial neuropsychological test scores, who have not undergone biofluid or imaging assessment for ADRD-related or other dementia-related pathology. Notably, from pharmacoeconomic, practical, and psychosocial standpoints, use of ZCoR for "long-range" clinical prognostication may be compatible with the up-to-decades-long, pre-clinical progression of beta-amyloid and tau neuropathology in AD: even 10 years before overt cognitive impairment, biofluid testing or imaging performed due to ZCoR high-risk status is likely to be informative[7], and the ZCoR classification, actionable. A second potential ZCoR application could be screening for undiagnosed prevalent cases of ADRD in primary care settings. Considering the estimated 45%–80% of dementia cases in older adults that go undiagnosed in the US[54], availability of a non-invasive, inexpensive, near-instantaneous, and almost universally-accessible tool could revolutionize detection of such patients. Third, ZCoR could be applied in scientific research regarding ADRD natural history and prevention. Beyond enrichment of trials of prophylactic interventions against cognitive impairment, ZCoR opens intriguing avenues of investigation, *e.g.*, examination of the roles of previously-underrecognized comorbidity classes with important associations with ADRD, *e.g.*, musculoskeletal disorders in males, respiratory infections in females, reproductive or ophthalmological disorders in both sexes. More precise understanding of the particular diseases that indeed are associated with ADRD will facilitate assessment of intriguing hypotheses such as inflammation serving as a key link between comorbidities and ADRD genetic features and phenotype[55].

## Limitations & Conclusion

Our key limitations arise from potential diagnostic mis-codings, and the current imprecision in Alzheimer's-related nomenclature[7]. Coupled with the high prevalence of undiagnosed dementia, mis-coding could lead to our ADRD signature deriving from data of only a fraction, albeit a substantial fraction, of our true cases. This situation might pose a particular peril under our Dx target definition, which considers only diagnostic codes. Mitigating this concern is the vast size of our control groups (n=375,101 females, n=330,921 males), implying that "non-signal" from large numbers of "true controls" is likely to overwhelm "buried ADRD signal" from "false controls". An additional possible concern related to mis-coding would be inclusion of non-ADRD age-related dementia cases among the ADRD group. However, given the "mixed" picture of dementia afflicting many patients with ADRD[1,7], tracking the characteristics of patients with non-ADRD cognitive impairment is also pertinent.

The performance of ZCoR might be further enhanced with the inclusion of treatment-related factors, *e.g.*, medications, along with comorbidities. As noted, however, using only diagnostic codes may increase the availability of data inputs for ZCoR in everyday practice, and hence the tool's scalability to routine settings. Moreover, comorbidity codes may be viewed to at least some extent as surrogates capturing the effects of medications that might influence Alzheimer neuropathology, *e.g.*, statins or anti-diabetic agents.

Predictive screening for ADRD raises some ethical concerns. In particular, early detection of progressive, not-yet-well-manageable brain disorders that have major effects on capacity, autonomy, and healthcare and other resource utilization, poses potential risks stemming from the possibility of stigmatization and discrimination[8]. It will be necessary to further explore these and other potential harms of early recognition of Alzheimer cognitive

impairment, and to seek their amelioration through legal and public health policy changes[3,8].

It is important to note that however strong its predictive performance, ZCoR is a screening tool, not a diagnostic tool, and and by itself certainly does not establish an ADRD diagnosis. ZCoR prediction of high-risk for ADRD should lead to diagnostic testing, $e.g.$, application of cognitive tests or imaging, and to intensified surveillance, as indicated. ZCoR results also can inform discussion with, and planning by, patients and their significant others.

In conclusion, ZCoR opens potentially new avenues in identification of and intervention against cognitive impairment, in neurocognitive research, and in designing effective caregiver support. Moving forward, we will focus on: 1) prospective validation of ZCoR; 2) assessment of the effects of ZCoR use on patient and caregiver quality-of-life, patient cognition and function, and healthcare utilization; 3) correlation with ADRD clinical and neuropathological biomarkers such as neuropsychological and functional test results and biofluid and imaging findings related to beta-amyloid, tau, and neurodegeneration; 4) comparison of ZCoR prospective performance in different racial groups and ethnicities, including examination of the signature's ability to reduce disparities in the rate of diagnosis. The impact of ZCoR on the accuracy and speed of diagnosis, on health care resource utilization, and eventually, on patient and caregiver outcomes, warrant prospective study.

## ACKNOWLEDGEMENTS

## AUTHOR CONTRIBUTIONS

DO implemented the algorithm and ran validation tests. DO and IC carried out mathematical modeling, and algorithm design. DO, SS, KR, JM and IC interpreted results. JM and IC guided the research. DO, KR, JM and IC wrote the paper. IC procured funding for the study.

## DECLARATION OF INTERESTS

IC is a founder and shareholder of Zero Burden Laboratories, Inc. He has not drawn any salary from the company. IC has received funding from the Alzheimer's Association, United States Department of Defense, the National Institutes of Health, and the Neubauer Collegium for Culture and Society. DO is a founder and shareholder of Zero Burden Laboratories. He has not drawn any salary from the company.

## STAR METHODS

### Resource Availability

*Materials availability*

This study did not generate or use new unique reagents.

*Data and code availability*

- **Data:** The Truven database used in this study is not in the public domain, and may be procured under license from https://www.ibm.com/watson-health/about/truven-health-analytics. A small de-identified set of patient diagnostic history is made available for testing purposes, and is publicly available as of the date of publication, as noted in the Key Resources Table. Description of ICD codes are available at https://www.cdc.gov/nchs/icd/icd10cm.htm, and https://www.icd10data.com/. A comprehensive interface for looking up ICD-10-CM code descriptions is provided by the National Library of Medicine, and may be accessed at https://clinicaltables.nlm.nih.gov/apidoc/icd10cm/v3/doc.html.
- **Code:** Working modules have been deposited at Zenodo and is publicly available as of the date of publication. DOIs are listed in the key resources table. Complete pseudocode is made available in the Supplementary Information text (Algorithms 1, and 2 in Supplementary Information).
- Any additional information required to reanalyze the data reported in this paper is available from the lead contact upon request.

*Lead Contact*

Further information and requests for resources and software should be directed to and will be fulfilled by the lead contact, Ishanu Chattopadhyay (ishanu@uchicago.edu).

## Method Details

### *Summary of Modeling Steps*

Individual diagnostic histories can have long-term memory[56], implying that the order, frequency, and comorbid interactions between diseases are important for assessing the future risk of our target phenotype. The large number of possible ICD codes, along with the sparsity of codes per patient (approximately one entry every 100 steps on the diagnostic time series) makes this a difficult learning problem.

### *Step 1: Partitioning The Human Disease Spectrum*

We begin by partitioning the human disease spectrum into 45 non-overlapping categories. Each category is defined by a set of diagnostic codes from the International Classification of Diseases, Ninth Revision (ICD9) (See Table SI-II for description of the categories used in this study).

For this study, we ended up using 6462501 and 9426722 diagnostic codes for males and females respectively (17501 and 18633 unique codes) spanning both ICD9 and ICD10 protocols (using ICD10 General Equivalence Mappings (GEMS)[57] equivalents where necessary), from a total 729,018 patients. Transforming the diagnostic histories to report only the broad categories reduces the number of distinct codes that the pipeline needs to handle, thus improving statistical power.

Our categories largely align with the top-level ICD9 categories, with small adjustments, *e.g.* bringing all infections under one category irrespective of the pathogen or the target organ. We do not pre-select the phenotypes; we want our algorithm to seek out the important patterns without any manual curation of the input data.

For each patient, the past medical history is a sequence $(t_1, x_1), \cdots, (t_m, x_m)$, where $t_i$ are timestamps and $x_i$ are ICD9 codes diagnosed at time $t_i$. We map individual patient history to a three-alphabet categorical time series $z^k$ corresponding to the disease category $k$, as follows. For each week $i$, we have:

$$z_i^k = \begin{cases} 0 & \text{if no diagnosis codes in week } i \\ 1 & \text{if there exists a diagnosis of category } k \text{ in week } i \\ 2 & \text{otherwise} \end{cases} \tag{1}$$

The time-series $z^k$ is observed in the inference period. Thus, each patient is represented by 43 mapped trinary series.

We refer to these individual diagnostic categories as "phenotypes", since they are observable characteristics of the patients. Each patient is represented by 45 sparse stochastic time series of events, which are compressed into specialized Hidden Markov Models known as Probabilistic Finite Automata (PFSA)[58,59]. These models are inferred separately for each phenotype, for each sex, and for the control and the positive cohorts, to identify the distinctive average patterns emerging at the population level. We infer $45 \times 2 \times 2 = 180$ PFSA models in total in this study. Our inference algorithm for these models does not presuppose a fixed structure, and is able to work with non-synchronized and variable-length data streams. Variation of these inferred models between the positive and control groups delineate the estimated risk of an ADRD diagnosis at the population level. Given these models, and the history of a specific patient, we can then quantify the likelihood of this patient's particular history being generated by the control PFSA models as opposed to the positive models. We refer to this likelihood difference as the sequence likelihood defect (SLD)[60], which is the one of the key informative features in our approach. The SLD is a novel concept, involving the generalization of the notion of Kullback-Liebler divergence[61] between probability distributions to a generalized divergence between possibly non-iid stochastic processes (See Step 2 below). SLD-based features allow the ZCoR measure to factor in complex longitudinal, *i.e.*, temporal patterns beyond simply the presence/absence of comorbidities.

### *Inference & Event Periods*

We train our predictive pipeline with all diagnostic codes that are recorded in the past 2 years from the point at which a prediction is made. This period from which we use data to train our pipeline is called the "inference window". Our aim is to make predictions on the occurrence of the target diagnostic codes at 1year from the end of the inference window. For patients in the control cohort, we make sure that no target code appears for 2years after the end of the inference window. Additionally, when making predictions further into the future, we always make sure that the control group has no target codes for 1 year after the predcited time of diagnosis,

*i.e.*, if we are making a prediction of a diagnosis $m$ years in future, then control group patients are chosen to have no diagnosis in at least next $m + 1$ years.

*Step 2: Model Inference & The Sequence Likelihood Defect $\Delta$*

The mapped series, disease-category, and ADRD diagnosis-status are considered to be independent sample paths, and we want to explicitly model these systems as specialized HMMs (PFSAs). We model the positive and the control cohorts and each disease category separately, ending up with a total of 86 HMMs at the population level (43 categories, 2 ADRD status categories: positive and control). Each of these inferred models is a PFSA; a directed graph with probability-weighted edges, and acts as an optimal generator of the stochastic process driving the sequential appearance of the three letters (as defined by Eq. (1)) corresponding to disease category, and ADRD status-type (See **"Probabilsitic Finite State Automata Inference"** for background on PFSA inference).

To reliably infer the ADRD status-type of a new patient, *i.e*, the likelihood of a diagnostic sequence being generated by the corresponding ADRD status-type model, we generalize the notion of Kullbeck-Leibler (KL) divergence[61,62] between probability distributions to a divergence $\mathcal{D}_{\mathsf{KL}}(G\|H)$ between ergodic stationary categorical stochastic processes[63] $G, H$ as:

$$\mathcal{D}_{\mathsf{KL}}(G\|H) = \lim_{n \to \infty} \frac{1}{n} \sum_{x:|x|=n} p_G(x) \log \frac{p_G(x)}{p_H(x)} \tag{2}$$

where $|x|$ is the sequence length, and $p_G(x), p_H(x)$ are the probabilities of sequence $x$ being generated by the processes $G, H$ respectively. Defining the log-likelihood of $x$ being generated by a process $G$ as :

$$L(x, G) = -\frac{1}{|x|} \log p_G(x) \tag{3}$$

The cohort-type for an observed sequence $x$ — which is actually generated by the hidden process $G$ — can be formally inferred from observations based on the following provable relationships (See Theorems 1 and 2):

$$\lim_{|x| \to \infty} L(x, G) = \mathcal{H}(G) \tag{4a}$$

$$\lim_{|x| \to \infty} L(x, H) = \mathcal{H}(G) + \mathcal{D}_{\mathsf{KL}}(G\|H) \tag{4b}$$

where $\mathcal{H}(\cdot)$ is the entropy rate of a process[61]. Importantly, Eq. (4) shows that the computed likelihood has an additional non-negative contribution from the divergence term when we choose the incorrect generative process. Thus, if a patient is eventually going to be diagnosed with ADRD, then we expect that the disease-specific mapped series corresponding to her diagnostic history be modeled by the PFSA in the positive cohort. Denoting the PFSA corresponding to disease category $j$ for positive and control cohorts as $G_+^j, G_0^j$ respectively, we can compute the *sequence likelihood defect* (SLD, $\Delta^j$) as:

$$\Delta^j \triangleq L(G_0^j, x) - L(G_+^j, x) \to \mathcal{D}_{\mathsf{KL}}(G_0^j\|G_+^j) \tag{5}$$

With the inferred PFSA models and the individual diagnostic history, we estimate the SLD measure on the right-hand side of Eqn. (5). The higher this likelihood defect, the higher the similarity of diagnosis history to that of women with ADRD.

*Step 3: Risk Estimation Pipeline With Semi-supervised & Supervised Learning Modules*

The risk estimation pipeline operates on patient specific information limited to the available diagnostic history in the inference period, and produces an estimate of the relative risk of ADRD, with an associated confidence value. To learn the parameters and associated model structures of this pipeline, we transform the patient specific data to a set of engineered features, and the feature vectors realized on the positive and control sets are used to train a gradient-boosting classifier[64]. The complete list of 701 features used is provided in Tab. VI.

We need two training sets: one to infer the models, and one to train the classifier with features derived from the inferred models. Thus, we do a random 3-way split of the set of unique patients into *feature-engineering* (25%), *training* (25%) and *test* (50%) sets. We use the feature-engineering set of ids first to infer our PFSA models *(unsupervised model inference in each category)*, which then allows us to train the gradient-boosting classifier using the training set and PFSA models *(classical supervised learning)*, and we finally execute out-of-sample validation on the test set. Fig. 2c in the main text shows the top 20 features ranked in order of their relative importance (relative loss in performance when dropped out of the analysis).

In addition to the phenotype specific specialized Markov models, we use a range of engineered features reflecting various aspects of diagnostic histories:

*Prevalence scores (p-scores):* The p-scores focus on individual diagnostic codes, and we create a dictionary of the ratio of relative prevalence of each code (relative to the set of all codes present) in the positive category (for each sex) to the control category. This is the second hyper-training step. In the later steps of the pipeline,

we use dictionary look ups to map codes to their p-scores, and also their aggregate measures such as mean, median, and variance to train a downstream LGBM.

*Rare scores:* These scores consist of a subset of p-scores which correspond to codes with particularly high and low relative prevalences (p-score $> 2$ or $< .5$). Thus, this feature category depends on the p-score dictionary generated in the second hyper-training step.

*Sequence scores:* Sequence scores are relatively straight-forward statsitical measures such as mean, median, variance, time since last occurrence $etc..$ on the trinary phenotype-specific sex-stratified histories. No hyper-training is required for the generation of the sequence features.

Thus we require three splits of the training dataset. The first split is used to carry out hyper-training of the PFSA models and the p-score dictionary. The second split is used to train the score-category specific LGBMs, one for each feature category. And the third split is used to train the final LGBM that takes inputs from the outputs of the four LGBMs in the previous layer. The network layout is shown in SI-Fig. 2.

### Validation

In validation, or actual prediction of patient fate, we use the trinary mapping, generate the features using the PFSA models and the p-score dictionary, and calculate the raw-risk via the trained LGBM network. The relative score is then obtained by a choice of the operating point reflecting the specificity/sensitivity trade-off discussed before.

### Data Splits: Training and Validation Hold-out

All eligible patients are randomly split into the Training set ($\approx 75\%$ of data) and the Test set ($\approx 25\%$ of data). The Training set is then split into 3 subsets: 1) The hyper-training set (SI-Fig. 2A) is used to train PFSA models p-score dictionary, 2) the second split (referred to as the pre-aggregation split, SI-Fig. 2B) is used to train the four feature-category specific LGBMs, and 3) the final split (referred to as the aggregation split, SI-Fig. 2C) is used to train the aggregate LGBM which uses outputs from the trained LGBMs in the previous layer. This trained pipeline is then validated on the held out validation split ($\approx 25\%$ of data).

### Generating PFSA Models From Set of Input Streams with Variable Input Lengths

Our PFSA reconstruction algorithm[59] is distinct from standard HMM learning. We do not need to pre-specify structures, or the number of states in the algorithm, and all model parameters are inferred directly from data. Additionally, we can operate either with 1) a single input stream, or 2) a set of input streams of possibly varying lengths which are assumed to be different and independent sample paths from the unknown stochastic generator we are trying to infer. At an intuitive level, we use the input data to infer the length of histories one must remember to estimate the current state, and predict futures for the process being modeled. Thus, we do not step through the symbol streams with a pre-specified model structure, and avoid the need to have equal-length inputs. More details of the algorithm are provided in the next section.

The ability to model a set of input streams of varying lengths is particularly important, since medical histories of different patients are typically of different lengths.

### Probabilsitic Finite State Automata Inference

Software for PFSA inference is mader available at https://pypi.org/project/zedsuite/.

Let $\Sigma$ be a finite alphabet of symbols with size $|\Sigma|$. The set of sequences of length $d$ over $\Sigma$ is denoted by $\Sigma^d$. The set of finite but unbounded sequences over $\Sigma$ is denoted by $\Sigma^\star$, the Kleene star operation[65], $i.e.$ $\Sigma^\star = \bigcup_{d=0}^\infty \Sigma^d$. We use lower case Greek, for example $\sigma$ or $\tau$, for symbols in $\Sigma$, and lower case Latin, for example $x$ or $y$, for sequences of symbols, $i.e.$ $x = \sigma_1 \sigma_2 \ldots \sigma_n$. We use $|x|$ to denote the length of $x$. The empty sequence is denoted by $\lambda$.

We denote the set of strictly infinite sequences over $\Sigma$ by $\Sigma^\omega$, and the set of strictly infinite sequences having $x$ as prefix by $x\Sigma^\omega$. Let $\mathcal{S} = \{x\Sigma^\omega : x \in \Sigma^\star\} \cup \{\emptyset\}$, we can verify that $\mathcal{S}$ is a semiring[66] over $\Sigma^\omega$. We use $\mathcal{F}$ to denote the sigma algebra generated by $\mathcal{S}$.

**Definition 1** (Stochastic Process over $\Sigma$). *A stochastic process over a finite alphabet $\Sigma$ is a collection of $\Sigma$-valued random variables $\{X_t\}_{t \in \mathbb{N}}$ indexed by positive integers[67].*

We are specifically interested in processes in which the $X_i$s are not necessarily independently distributed.

**Definition 2** (Sequence-Induced Measure and Derivative). *For a process $\mathscr{P}$, let $\Pr_{\mathscr{P}}(x)$ or simply $\Pr(x)$ denote the probability $\mathscr{P}$ producing a sample path prefixed by $x$. The **measure** $\mu_x$ **induced by a sequence** $x \in \Sigma^\star$ is the extension[66] to $\mathcal{F}$ of the premeasure defined on the semiring $\mathbb{S}$ given by*

$$\forall x, y \in \Sigma^\star, \mu_x(y\Sigma^\omega) \triangleq \frac{\Pr(xy)}{\Pr(x)}, \text{ if } \Pr(x) > 0 \tag{6}$$

*For any $d \in \mathbb{N}$, the $d$-**th order derivative** of a sequence $x$, written as $\phi_x^d$, is defined to be the marginal distribution of $\mu_x$ on $\Sigma^d$, with the entry indexed by $y$ denoted by $\phi_x^d(y)$. The first-order derivative is called the **symbolic derivative** and is denoted by $\phi_x$ for short.*

**Definition 3** (Probabilistic Nerode Equivalence and Causal States[68]). *For any pair of sequences $x, y \in \Sigma^\star$, $x$ is equivalent to $y$, written as $x \sim y$, if and only if either $\Pr(x) = \Pr(y) = 0$, or $\mu_x = \mu_y$. The equivalence class of a sequence $x$ is denoted by $[x]$ and is called a **causal state**[69]. The cardinality of the set of causal states is called the **probabilistic Nerode index**, or the Nerode index for simplicity.*

We can see from the definition that causal states captures how the history of a process influences its future. Since the probabilistic Nerode equivalence is right invariant, it gives rise naturally to a automaton structure introduced below.

**Definition 4** (Probabilistic Finite-State Automaton (PFSA)). *A PFSA $G$ is defined by a quadruple $(Q, \Sigma, \delta, \widetilde{\pi})$, where $Q$ is a finite set, $\Sigma$ is a finite alphabet, $\delta : Q \times \Sigma \to \Sigma$ is called the transition map, and $\widetilde{\pi} : Q \to \mathbf{P}_\Sigma$, where $\mathbf{P}_\Sigma$ is the space of probability distributions over $\Sigma$, is called the transition probability. The entry of $\widetilde{\pi}(q)$ indexed by $\sigma$ is denoted by $\widetilde{\pi}(q, \sigma)$.*

**Definition 5** (Transition and Observation Matrices). *The transition matrix $\Pi$ is the $|Q| \times |Q|$ matrix with the entry indexed by $q, q'$, written as $\pi_{q,q'}$, satisfying*

$$\pi_{q,q'} \triangleq \sum_{\{\sigma \in \Sigma | \delta(q,\sigma)=q'\}} \widetilde{\pi}(q, \sigma) \tag{7}$$

*and the observation matrix $\widetilde{\Pi}$ is a $|Q| \times |\Sigma|$ matrix with the entry indexed by $q, \sigma$ equaling $\widetilde{\pi}(q, \sigma)$.*

We note that both $\Pi$ and $\widetilde{\Pi}$ are stochastic, *i.e.* non-negative with rows summing up to 1.

**Definition 6** (Extension of $\delta$ and $\widetilde{\pi}$ to $\Sigma^\star$). *For any $x = \sigma_1 \ldots \sigma_k$, $\delta(q, x)$ is defined recursively by*

$$\delta(q, x) \triangleq \delta(\delta(q, \sigma_1 \ldots \sigma_{k-1}), \sigma_k) \tag{8}$$

*with $\delta(q, \lambda) = q$, and $\widetilde{\pi}(q, x)$ is defined recursively by*

$$\widetilde{\pi}(q, x) \triangleq \prod_{i=1}^k \widetilde{\pi}(\delta(q, \sigma_1 \ldots \sigma_{i-1}), \sigma_i) \tag{9}$$

*with $\widetilde{\pi}(q, \lambda) = 1$.*

**Definition 7** (Strongly Connected PFSA). *We say a PFSA is strongly connected if the underlying directed graph is strongly connected[70]. More precisely, a PFSA $G = (Q, \Sigma, \delta, \widetilde{\pi})$ is strongly connected if for any pair of distinct states $q$ and $q' \in Q$, there is an $x \in \Sigma^\star$ such that $\delta(q, x) = q'$.*

We assume all PFSA in the discussions in the sequel are strongly connected if not specified otherwise. For strongly connected PFSA $G$, there is a unique probability distribution over $Q$ that satisfies $\mathbf{v}^\mathsf{T}\Pi = \mathbf{v}^\mathsf{T}$. This is the **stationary distribution**[71,72] of $G$ and is denoted as $\wp_G$, or $\wp$ if $G$ is understood.

**Definition 8** (Γ-Expression). *We can encode the information contained in $\delta$ and $\widetilde{\pi}$ by a set of $|Q| \times |Q|$ matrices $\Gamma = \{\Gamma_\sigma | \sigma \in \Sigma\}$, where*

$$\Gamma_\sigma\big|_{q,q'} \triangleq \begin{cases} \widetilde{\pi}(q, \sigma) & \text{if } \delta(q, \sigma) = q', \\ 0 & \text{if otherwise.} \end{cases} \tag{10}$$

*$\Gamma_\sigma$ is called **event-specific transition matrix**, with the event being that $\sigma$ is current the output. $\Gamma_\sigma$ can also be extended to arbitrary $x \in \Sigma^\star$ by defining $\Gamma_x = \prod_{i=1}^k \Gamma_{\sigma_i}$ with $\Gamma_\lambda = I$.*

**Definition 9** (Sequence-Induced Distribution on States). *For a PFSA $G = (Q, \Sigma, \delta, \widetilde{\pi})$ and a distribution $\wp_0$ on $Q$, the **distribution on $Q$ induced by a sequence** $x$ is given by $\wp_{G,\wp_0}^\mathsf{T}(x) = \llbracket \wp_0^\mathsf{T} \Gamma_x \rrbracket$ with $\wp_{G,\wp_0}(\lambda) = \wp_0$. The entry indexed by $q \in Q$ of the vector $\wp_{G,\wp_0}(x)$ is written as $\wp_{G,\wp_0}(x, q)$. When $\wp_0 = \wp_G$, the stationary distribution of $G$, we write $\wp_{G,\wp_0}(x)$ as $\wp_G(x)$, or simply as $\wp(x)$, if $G$ is understood.*

**Definition 10** (Stochastic Process Generated by a PFSA). *Let $G = (Q, \Sigma, \delta, \widetilde{\pi})$ be a PFSA and let $\wp_0$ be a distribution on $Q$, the $\Sigma$-valued stochastic process $\{X_t\}_{t \in \Sigma}$ generated by $G$ and $\wp_0$ satisfies that $X_1$ follows the distribution $\wp_0$ and $X_{t+1}$ follows the distribution $\wp_{G,\wp_0}(X_1 \cdots X_t)$ for $t \in \mathbb{N}$.*

For the rest of this paper, we will assume $\wp_0 = \wp_G$ if not specified otherwise. We can show that, when initialized with $\wp_G$, the process generated by a PFSA G is stationary and ergodic. We also note the, for the process generate by G, we have $\phi_x = \wp_G(x)^\mathsf{T}\widetilde{\Pi}$. Since $\wp_G(\lambda) = \wp_G$, the symbolic derivative of the empty sequence $\phi_\lambda$ is the stationary distribution on the symbols.

**Definition 11** (Synchronizable PFSA and Synchronizing Sequence). *A **synchronizing sequence** is a finite sequence that sends an arbitrary state of the PFSA to a fixed state[73]. To be more precise, let $G = (Q, \Sigma, \delta, \widetilde{\pi})$ be a PFSA, we say a sequence $x \in \Sigma^\star$ is a synchronizing sequence to a state $q \in Q$ if $\delta(q', x) = q$ for all $q' \in Q$. A PFSA is **synchronizable** if it has at least one synchronizing sequence. Given a sample path generated by a PFSA, we say the PFSA is **synchronized** if a synchronizing sequence transpires in the sample path.*

**Definition 12** (Equivalence and Irreducibility). *Two PFSA G and H are **equivalent** if they generate the same stochastic process. A PFSA G is said to be **irreducible**, if there is not another PFSA with smaller state set that is equivalent to G.*

**Definition 13.** *Consider a PFSA G over state set Q. For a give $\varepsilon > 0$, we say a sequence $x$ is a $\varepsilon$-synchronizing sequence to a state $q \in Q$ if*

$$\|\wp_G(x) - \mathbf{e}_q\|_\infty \leqslant \varepsilon. \tag{11}$$

While there exists PFSA that is not synchronizable, we can show that an irreducible PFSA always has an $\varepsilon$-synchronizing sequence for some state q for arbitrarily small $\varepsilon > 0$. Moreover, we can show that as length increases, sequences produced by PFSA become uniformly $\varepsilon$-synchronizing. These two are the underpinning properties for the inference algorithm of PFSA (See Alg. 1), because they imply that $\phi_x$ can be used to approximate $\widetilde{\pi}(q)$ if $x$ are properly prefixed and long enough.

**Definition 14** (Joint $\varepsilon$-Synchronizing Sequence). *Let G and H be two PFSA over state sets $Q_G$ and $Q_H$, respectively. For a fixed $\varepsilon$, a sequence $x$ is said to be **jointly $\varepsilon$-synchronizing** to $(q, r) \in Q_G \times Q_H$ if $x$ is $\varepsilon$-synchronizing to q and to r simultaneously. We define*

$$\Sigma_{\varepsilon,(q,r)}^d \triangleq \left\{ x \in \Sigma^d : x \text{ jointly } \varepsilon\text{-synchronizing to } (q,r) \right\} \tag{12}$$

**Definition 15** (Joint Pair of States). *Let G and H be two PFSA over state sets $Q_G$ and $Q_H$, respectively. Define*

$$p_G(q, r) \triangleq \lim_{d \to \infty} p_G\left( \Sigma_{\varepsilon,(q,r)}^d \right) \tag{13}$$

*A pair of states $(q, r) \in Q_G \times Q_H$ is called a G-**joint pair** of states if $p_G(q, r) > 0$. We also define*

$$Q_c \triangleq \{(q, r) \in Q_G \times Q_H : (q, r) \text{ is a G-joint pair}\} \tag{14}$$

The inference algorithm for PFSA is called `GenESeSS` for <u>Gen</u>erator <u>E</u>xtraction Using <u>Se</u>lf-<u>s</u>imilar <u>S</u>emantics. With an input sequence $x$ and a hyperparameter $\varepsilon$, `GenESeSS` outputs a PFSA in the following three steps: 1) approximate an almost synchronizing sequence; 2) identify the transition structure of the PFSA; 3) calculate the transition probabilities of the PFSA. See Alg. 1[59] for details.

*Theoretical Development of Sequence Likelihood Defect*

**Definition 16** (Entropy Rate and KL Divergence). *By entropy rate of a PFSA, we mean the entropy rate of the stochastic process generated by the PFSA[74]. Similarly, by KL divergence of two PFSA, we mean the KL divergence between the two processes generated by them[75]. More precisely, we have*

$$\mathcal{H}(G) = -\lim_{d \to \infty} \frac{1}{d} \sum_{x \in \Sigma^d} p(x) \log p(x) \tag{15}$$

*and the KL divergence*

$$\mathcal{D}_{KL}(G \| H) = \lim_{d \to \infty} \frac{1}{d} \sum_{x \in \Sigma^d} p_G(x) \log \frac{p_G(x)}{p_H(x)} \tag{16}$$

*whenever the limits exist.*

**Theorem 1** (Closed-form Formula for Entropy Rate and KL Divergence). *The entropy rate of a PFSA $G = (\Sigma, Q, \delta, \widetilde{\pi})$ is given by*

$$\mathcal{H}(G) = \sum_{q \in Q} \wp_G(q) \cdot h(\widetilde{\pi}(q)) \tag{17}$$

*where $h(\mathbf{v})$ is the based-2 entropy of the probability vector $\mathbf{v}$.*

*Consider two PFSA $G = (Q_G, \Sigma, \delta_G, \widetilde{\pi}_G)$ and $H = (Q_H, \Sigma, \delta_H, \widetilde{\pi}_H)$ with $\mu_G$ being absolutely continuous with*

respect to $\mu_H$. Let $Q_c$ be the set of $G$-joint pairs of states, we have

$$\mathcal{D}_{KL}(G \| H) = \sum_{(q,r) \in Q_c} p_G(q,r) D_{KL}(\widetilde{\pi}_G(q) \| \widetilde{\pi}_H(r))$$

(18)

**Definition 17** (Log-likelihood). *Let* $x \in \Sigma^d$, *the log-likelihood[74] of a PFSA* $G$ *generating* $x$ *is given by*

$$L(x, G) = -\frac{1}{d} \log p_G(x)$$

(19)

The calculation of log-likelihood is detailed in Alg. 2.

**Theorem 2** (Convergence of log-likelihood). *Let* $G$ *and* $H$ *be two reduced PFSA, and let* $x \in \Sigma^d$ *be a sequence generated by* $G$. *Then we have*

$$L(x, H) \to \mathcal{H}(G) + \mathcal{D}_{KL}(G \| H)$$

(20)

*in probability as* $d \to \infty$.

*Proof.* We first notice that

$$\sum_{x \in \Sigma^d} p_G(x) \log \frac{p_G(x)}{p_H(x)} = \sum_{x \in \Sigma^{d-1}} \sum_{\sigma \in \Sigma} p_G(x) \wp_G(x) \widetilde{\Pi}_G \Big|_\sigma \log \frac{p_G(x) \wp_G(x) \widetilde{\Pi}_G \Big|_\sigma}{p_H(x) \wp_H(x) \widetilde{\Pi}_H \Big|_\sigma}$$

(21)

$$= \sum_{x \in \Sigma^{d-1}} p_G(x) \log \frac{p_G(x)}{p_H(x)} + \underbrace{\sum_{x \in \Sigma^{d-1}} p_G(x) \sum_{\sigma \in \Sigma} \wp_G(x) \widetilde{\Pi}_G \Big|_\sigma \log \frac{\wp_G(x) \widetilde{\Pi}_G \Big|_\sigma}{\wp_H(x) \widetilde{\Pi}_H \Big|_\sigma}}_{D_d}$$

(22)

By induction, we have $\mathcal{D}_{KL}(G \| H) = \lim_{d \to \infty} \frac{1}{d} \sum_{i=1}^d D_i$, and hence by Cesàro summation theorem[76], we have $\mathcal{D}_{KL}(G \| H) = \lim_{d \to \infty} D_d$. Let $x = \sigma_1 \sigma_2 ... \sigma_n$ be a sequence generated by $G$. Let $x^{[i-1]}$ is the truncation of $x$ at the $(i-1)$-th symbols, we have

$$-\frac{1}{n} \sum_{i=1}^n \log \wp_H\left(x^{[i-1]}\right) \widetilde{\Pi}_H \Big|_{\sigma_i} = \underbrace{\frac{1}{n} \sum_{i=1}^n \log \frac{\wp_G\left(x^{[i-1]}\right) \widetilde{\Pi}_G \Big|_{\sigma_i}}{\wp_H\left(x^{[i-1]}\right) \widetilde{\Pi}_H \Big|_{\sigma_i}}}_{A_{x,n}} - \underbrace{\frac{1}{n} \sum_{i=1}^n \log \wp_G\left(x^{[i-1]}\right) \widetilde{\Pi}_G \Big|_{\sigma_i}}_{B_{x,n}}$$

(23)

Since the stochastic process $G$ generates is ergodic, we have

$$\lim_{n \to \infty} A_{x,n} = \lim_{d \to \infty} D_d = \mathcal{D}_{KL}(G \| H)$$

(24)

and $\lim_{n \to \infty} B_{x,n} = \mathcal{H}(G)$. □

## Quantification & Statistical Analysis

### Raw Risk & Relative Risk

We choose a decision threshold on the raw risk computed by our pipeline to make crisp predictions, *i.e.*, if the raw risk is greater than this calibrated threshold, then the individual patient is predicted to be in the positive category.

### Threshold Selection on ROC Curve

In situations where the number of negatives vastly outnumber the number of positives (which is the case in our problem), it is better to base this trade-off on a measure that is independent of the number of true negatives. The two popular measures considered in the literature are accuracy and the F1-score:

$$\text{accuracy} = \frac{t_p + t_n}{t_p + f_p + f_n + t_n}$$

(25)

$$F1 = \frac{2 t_p}{2 t_p + f_p + f_n}$$

(26)

The F1-score is the same as accuracy where the number of true negatives is the same as the number of true positives, thus partially correcting for the class imbalance.

The selection of the threshold may also be dictated by the current practice of ensuring high specificities in screening tests. Thus, the most relevant clinically operating point is either the one corresponding to 95% specificity, which is highlighted in Fig. 2a.

*Performance Measurement*

We measure our performance using standard metrics including the AUC, sensitivity, specificity, the positive predictive value (PPV), and the negative predictive value (NPV). We also report accuracy (acc, See Tables VII and VIII), which is the probability of a correct prediction (positive or control), and variation of AUC for predicting ADRD into the future up to 10 years (See Table IX).

Ninety-five percent confidence intervals (95% CIs) on ROC curves and AUCs were obtained via bootstrapping, and AUC p-vales were obtained using the Mann-Whitney U-test statistic.

*Note on Reciever Operating Characteristics (ROC) and Precision-recall Curves*

The ROC curve is a plot between the False Positive rate (TPR) and the True Positive Rate (TPR), and the area under the ROC curve (AUC) is often used as a measure of classifier performance. For the same of completeness, we introduce the relevant definitions:

In the following $P$ denotes the total number of positive samples (number of patients who are eventually diagnosed), and $N$ denotes the total number of negative samples (number of patients in the control group).

**Definition 18.** *True positive rate, true negative rate, false positive rate, positive predictive value (**PPV**), and* **prevalence** *($\rho$) are defined as:*

$$\text{sensitivity, or } TPR = \frac{t_p}{P} = \frac{t_p}{t_p + f_n} \tag{27}$$

$$\text{specificity, or } TNR = \frac{t_n}{N} = \frac{t_n}{t_n + f_p} \tag{28}$$

$$FPR = 1 - TNR \tag{29}$$

$$\text{precision, or } PPV = \frac{t_p}{t_p + f_p} \tag{30}$$

$$\rho = \frac{P}{N + P} \tag{31}$$

*where as before* $t_p, t_n, f_p, f_n$ *are true positives, true negatives, false positives, and false negatives respectively.*

Denoting sensitivity by $s$, and specifciyty by $c$, it follows that:

$$PPV = \frac{t_p/P}{t_p/P + (f_p/N)(N/P)} = \frac{TPR}{TPR + ((N - t_n)/N)(N/P)} \tag{32}$$

$$\Rightarrow PPV = \frac{s}{s + (1 - c)(\frac{1}{\rho} - 1)} \tag{33}$$

Thus, we note that for a fixed specificity and sensitivity, the PPV depends on prevalence. Indeed, it is clear from the above argument that PPV decreases with decreasing prevalence, and vice versa.

*Effect of Class Imbalance*

ROC curves are generally assumed to be robust to class imbalance. Note that if we assume that patient outcomes are independent (which is well-justified in the case of a non-communicable condition, particularly in large databases), then $t_p$ should scale linearly with the total number of positives $P$, implying:

$$TPR = \frac{t_p}{P} = \frac{t'_p}{P'} \tag{34}$$

implying that with different sizes of the set of positive samples (or negative samples), the ROC curve remains unchanged. In particular, note that even if the prevalence is very small (say 0.01%), we cannot cheat to boost the AUC by labeling all predictions as negative, or stating that risk is always zero: in that case, our $P$ is very small, but our $t_p = 0$ strictly, implying that our $TPR = 0$, thus leading to a zero AUC. We can cheat to boost the accuracy (See the previous section), but not the AUC.

Note that while relative class sizes or imbalance does not affect the ROC (under the assumption that true positives and true negatives scale with the number of positives and negatives), very small absolute sample sizes might still result in poor performance of the model.

The precision-recall curves do get affected by class imbalance, or the prevalence, as shown by Eq (33). However, in diagnostic analysis, they are important since we are generally less interested in the number of true negatives; the ratio of false positives to the total number of positive recommendations by the algorithm is much more relevant, *i.e.*, the PPV or the precision.

# REFERENCES

[1] Arvanitakis, Z., Shah, R. C. & Bennett, D. A. Diagnosis and management of dementia. *Jama* **322**, 1589–1599 (2019).

[2] Moyer, V. A. Screening for cognitive impairment in older adults: Us preventive services task force recommendation statement. *Annals of internal medicine* **160**, 791–797 (2014).

[3] Owens, D. K. *et al.* Screening for cognitive impairment in older adults: Us preventive services task force recommendation statement. *Jama* **323**, 757–763 (2020).

[4] Chu, L. *et al.* Alzheimer's disease: early diagnosis and treatment. *Hong Kong Med J* **18**, 228–237 (2012).

[5] Association, A. *et al.* Alzheimer's disease facts and figures. *Alzheimer's & dementia* **13**, 325–73 (2017).

[6] Murray, C. J. *et al.* The state of us health, 1990-2010: burden of diseases, injuries, and risk factors. *Jama* **310**, 591–606 (2013).

[7] Jack Jr, C. R. *et al.* NIA-AA research framework: toward a biological definition of Alzheimer's disease. *Alzheimer's & Dementia* **14**, 535–562 (2018).

[8] Ahlgrim, N. S., Garza, K., Hoffman, C. & Rommelfanger, K. S. Prodromes and preclinical detection of brain diseases: Surveying the ethical landscape of predicting brain health. *Eneuro* **6** (2019).

[9] Patnode, C. D. *et al.* Screening for cognitive impairment in older adults: updated evidence report and systematic review for the us preventive services task force. *Jama* **323**, 764–785 (2020).

[10] Borson, S. *et al.* Implementing routine cognitive screening of older adults in primary care: process and impact on physician behavior. *Journal of general internal medicine* **22**, 811–817 (2007).

[11] Davidson, M. & Thibaut, F. Is dementia a preventable disease? *Dialogues in clinical neuroscience* **21**, 3 (2019).

[12] Mattsson-Carlgren, N. *et al.* Longitudinal plasma p-tau217 is increased in early stages of alzheimer's disease. *Brain* **143**, 3234–3241 (2020).

[13] Moscoso, A. *et al.* Longitudinal associations of blood phosphorylated tau181 and neurofilament light chain with neurodegeneration in alzheimer disease. *JAMA Neurology* .

[14] Karikari, T. K. *et al.* Diagnostic performance and prediction of clinical progression of plasma phospho-tau181 in the Alzheimer's Disease Neuroimaging Initiative. *Molecular Psychiatry* 1–14 (2020).

[15] Hall, S. *et al.* Plasma phospho-tau identifies Alzheimer's co-pathology in patients with lewy body disease. *Movement Disorders* (2020).

[16] Janelidze, S. *et al.* Plasma p-tau181 in alzheimer's disease: relationship to other biomarkers, differential diagnosis, neuropathology and longitudinal progression to alzheimer's dementia. *Nature medicine* **26**, 379–386 (2020).

[17] Bezdicek, O. *et al.* Determining a short form montreal cognitive assessment (s-moca) czech version: validity in mild cognitive impairment parkinson's disease and cross-cultural comparison. *Assessment* **27**, 1960–1970 (2020).

[18] Lab, S. R. A. Montreal cognitive assessment in rehabmeasures database. https://www.sralab.org/rehabilitation-measures/montreal-cognitive-assessment (2020). (Accessed on 02/14/2021).

[19] Nasreddine, Z. S. *et al.* The Montreal Cognitive Assessment, MoCA: a brief screening tool for mild cognitive impairment. *Journal of the American Geriatrics Society* **53**, 695–699 (2005).

[20] Zhao, L. Alzheimer's disease facts and figures. *Alzheimers Dement* **16**, 391–460 (2020).

[21] Wilkinson, T. *et al.* Identifying dementia cases with routinely collected health data: a systematic review. *Alzheimer's & Dementia* **14**, 1038–1051 (2018).

[22] Duthie, A., Chew, D. & Soiza, R. Non-psychiatric comorbidity associated with Alzheimer's disease. *QJM: An International Journal of Medicine* **104**, 913–920 (2011).

[23] Anstey, K. J. *et al.* Future directions for dementia risk reduction and prevention research: An international research network on dementia prevention consensus. *Journal of Alzheimer's Disease* 1–10 (2020).

[24] Sipilä, P. N. *et al.* Hospital-treated infectious diseases and the risk of dementia: multicohort study with replication in the uk biobank. *medRxiv* (2020).

[25] Muzambi, R. *et al.* Assessment of common infections and incident dementia using UK primary and secondary care data: a historical cohort study. *The Lancet Healthy Longevity* (2021).

[26] Eavani, H. *et al.* Heterogeneity of structural and functional imaging patterns of advanced brain aging revealed via machine learning methods. *Neurobiology of aging* **71**, 41–50 (2018).

[27] Montreal Cognitive Assessment Webpage. Normative test. https://www.mocatest.org/moca-clinic-data/. (Accessed on 05/04/2021).

[28] Ferretti, M. T. *et al.* Sex and gender differences in alzheimer's disease: current challenges and implications for clinical practice: position paper of the dementia and cognitive disorders panel of the european academy of Neurology. *European journal of Neurology* **27**, 928–943 (2020).

[29] Choi, J., Kwon, L.-N., Lim, H. & Chun, H.-W. Gender-based analysis of risk factors for dementia using

senior cohort. *International journal of environmental research and public health* **17**, 7274 (2020).

[30] Exalto, L. G. *et al.* Risk score for prediction of 10 year dementia risk in individuals with type 2 diabetes: a cohort study. *The Lancet Diabetes & Endocrinology* **1**, 183–190 (2013).

[31] Reitz, C. *et al.* A summary risk score for the prediction of alzheimer disease in elderly persons. *Archives of Neurology* **67**, 835–841 (2010).

[32] Barnes, D. E. *et al.* Development and validation of a brief dementia screening indicator for primary care. *Alzheimer's & Dementia* **10**, 656–665 (2014).

[33] Tang, E. Y. *et al.* Current developments in dementia risk prediction modelling: an updated systematic review. *PloS one* **10**, e0136181 (2015).

[34] Chary, E. *et al.* Short-versus long-term prediction of dementia among subjects with low and high educational levels. *Alzheimer's & Dementia* **9**, 562–571 (2013).

[35] Ohara, T. *et al.* Apolipoprotein genotype for prediction of Alzheimer's disease in older japanese: the hisayama study. *Journal of the American Geriatrics Society* **59**, 1074–1079 (2011).

[36] Boustani, M. *et al.* Passive digital signature for early identification of Alzheimer's disease and related dementia. *Journal of the American Geriatrics Society* **68**, 511–518 (2020).

[37] Park, J. H. *et al.* Machine learning prediction of incidence of alzheimer's disease using large-scale administrative health data. *NPJ digital medicine* **3**, 1–7 (2020).

[38] So, A., Hooshyar, D., Park, K. W. & Lim, H. S. Early diagnosis of dementia from clinical data by machine learning techniques. *Applied Sciences* **7**, 651 (2017).

[39] Hansen, L. The truven health marketscan databases for life sciences researchers. *Truven Health Ananlytics IBM Watson Health* (2017).

[40] FDA. Fda-approved treatments for Alzheimer's disease or treatments for Alzheimer's disease. https://www.alz.org/media/documents/fda-approved-treatments-alzheimers-ts.pdf (2019). (Accessed on 02/14/2021).

[41] Tortajada-Soler, M. *et al.* Prevalence of comorbidities in individuals diagnosed and undiagnosed with alzheimer's disease in león, spain and a proposal for contingency procedures to follow in the case of emergencies involving people with alzheimer's disease. *International journal of environmental research and public health* **17**, 3398 (2020).

[42] Fouladvand, S. *et al.* Deep learning prediction of mild cognitive impairment using electronic health records. In *2019 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 799–806 (IEEE, 2019).

[43] Anheim, M. *et al.* Ataxic variant of alzheimer's disease caused by pro117ala psen1 mutation. *Journal of Neurology, Neurosurgery & Psychiatry* **78**, 1414–1415 (2007).

[44] Piccini, A. *et al.* Association of a presenilin 1 s170f mutation with a novel alzheimer disease molecular phenotype. *Archives of Neurology* **64**, 738–745 (2007).

[45] Humbert, I. A. *et al.* Early deficits in cortical control of swallowing in Alzheimer's disease. *Journal of Alzheimer's disease* **19**, 1185–1197 (2010).

[46] Kai, K. *et al.* Relationship between eating disturbance and dementia severity in patients with alzheimer's disease. *PloS one* **10**, e0133666 (2015).

[47] Gannon, O., Robison, L., Custozzo, A. & Zuloaga, K. Sex differences in risk factors for vascular contributions to cognitive impairment & dementia. *Neurochemistry international* **127**, 38–55 (2019).

[48] Kim, S. E. *et al.* Sex-specific relationship of cardiometabolic syndrome with lower cortical thickness. *Neurology* **93**, e1045–e1057 (2019).

[49] Elbejjani, M. *et al.* Depression, depressive symptoms, and rate of hippocampal atrophy in a longitudinal cohort of older men and women. *Psychological medicine* **45**, 1931 (2015).

[50] Hua, X. *et al.* Sex and age differences in atrophic rates: an adni study with n= 1368 mri scans. *Neurobiology of aging* **31**, 1463–1480 (2010).

[51] Irvine, K., Laws, K. R., Gale, T. M. & Kondel, T. K. Greater cognitive deterioration in women than men with Alzheimer's disease: a meta analysis. *Journal of clinical and experimental neuropsychology* **34**, 989–998 (2012).

[52] of Neurology, A. A. Practice guideline update summary: Mild cognitive impairment. https://www.aan.com/Guidelines/home/GuidelineDetail/881 (2018). (Accessed on 04/08/2021).

[53] Belleville, S., Fouquet, C., Hudon, C., Zomahoun, H. T. V. & Croteau, J. Neuropsychological measures that predict progression from mild cognitive impairment to Alzheimer's type dementia in older adults: a systematic review and meta-analysis. *Neuropsychology review* **27**, 328–353 (2017).

[54] Fowler, N. R. *et al.* Older primary care patients' attitudes and willingness to screen for dementia. *Journal of Aging Research* **2015** (2015).

[55] Newcombe, E. A. *et al.* Inflammation: the link between comorbidities, genetics, and Alzheimer's disease. *Journal of neuroinflammation* **15**, 1–26 (2018).

[56] Granger, C. W. J. & Joyeux, R. An introduction to long-memory time series models and fractional differencing. *Journal of Time Series Analysis* **1**, 15–29.

[57] General equivalence mappings (2019). URL https://www.cms.gov/Medicare/Coding/ICD10/downloads/ICD-10_GEM_fact_sheet.pdf.

[58] Chattopadhyay, I. & Ray, A. Structural transformations of probabilistic finite state machines. *International Journal of Control* **81**, 820–835 (2008).

[59] Chattopadhyay, I. & Lipson, H. Abductive learning of quantized stochastic processes with probabilistic finite automata. *Philos Trans A* **371**, 20110543 (2013).

[60] Huang, Y. & Chattopadhyay, I. Data smashing 2.0: Sequence likelihood (sl) divergence for fast time series comparison. *arXiv preprint arXiv:1909.12243* (2019).

[61] Cover, T. M. & Thomas, J. A. *Elements of Information Theory* (Wiley-Interscience, New York, NY, USA, 1991).

[62] Kullback, S. & Leibler, R. A. On information and sufficiency. *Ann. Math. Statist.* **22**, 79–86 (1951). URL https://doi.org/10.1214/aoms/1177729694.

[63] Doob, J. *Stochastic Processes*. Wiley Publications in Statistics (John Wiley & Sons, 1953). URL https://books.google.com/books?id=KvJQAAAAMAAJ.

[64] Friedman, J. H. Stochastic gradient boosting. *Comput. Stat. Data Anal.* **38**, 367–378 (2002). URL http://dx.doi.org/10.1016/S0167-9473(01)00065-2.

[65] Hopcroft, J. E. *Introduction to automata theory, languages, and computation* (Pearson Education India, 2008).

[66] Klenke, A. *Probability theory: a comprehensive course* (Springer Science & Business Media, 2013).

[67] Doob, J. *Stochastic processes*. Wiley publications in statistics (Wiley, 1990). URL https://books.google.com/books?id=7Bu8jgEACAAJ.

[68] Chattopadhyay, I. & Ray, A. Structural transformations of probabilistic finite state machines. *International Journal of Control* **81**, 820–835 (2008).

[69] Chattopadhyay, I. & Lipson, H. Data smashing: uncovering lurking order in data. *Journal of The Royal Society Interface* **11**, 20140826 (2014).

[70] Bondy, J. & Murty, U. Graph theory (2008). *Grad. Texts in Math* (2008).

[71] Vidyasagar, M. *Hidden markov processes: Theory and applications to biology*, vol. 44 (Princeton University Press, 2014).

[72] Kai, L. C. *Markov Chains: With Stationary Transition Probabilities* (Springer-Verlag, 1967).

[73] Trahtman, A. N. The road coloring and Černý conjecture. In *Proc. of Prague stringology conference*, vol. 1, 12 (Citeseer, 2008).

[74] Cover, T. M. & Thomas, J. A. *Elements of information theory* (John Wiley & Sons, 2012).

[75] Matthews, A. G. d. G., Hensman, J., Turner, R. & Ghahramani, Z. On sparse variational methods and the kullback-leibler divergence between stochastic processes. *Journal of Machine Learning Research* **51**, 231–239 (2016).

[76] Hardy, G. Divergent series, with a preface by je littlewood and a note by ls bosanquet, reprint of the revised (1963) edition. *Éditions Jacques Gabay, Sceaux* (1992).

# Key Resources Table

| REAGENT | SOURCE | IDENTIFIER |
|---|---|---|
| Deposited data | | |
| Truven Marketscan Database | IBM Watson® (with appropriate licensing) | https: //www.ibm.com/watson-health/about/truven-health-analytics |
| Small Patient Database | Excerpt from University of Chicago Medicine de-identified records between 2012-2021 | https://github.com/zeroknowledgediscovery/EHRdata https://doi.org/10.5281/.zenodo.5348229 |
| Software and algorithms | | |
| PFSA inference algorithm implementation | Laboratory of Zero Knowledge Discovery (zed.uchicago.edu) | https://pypi.org/project/zedsuite/ |
| ZCoR modules | Laboratory of Zero Knowledge Discovery (zed.uchicago.edu) | https://github.com/zeroknowledgediscovery/ZCOR-ADRD https://doi.org/10.5281/zenodo.5348219 |

# Supplementary Text: Rapid Universal Early Screening for Alzheimer's Disease and Related Dementia via Pattern Discovery in Diagnostic History

Dmytro Onishchenko[1], Sam Searle[7], Kenneth Rockwood[7], James A. Mastrianni[5,6] and Ishanu Chattopadhyay[1,2,3,4]★

[1]Department of Medicine, University of Chicago, Chicago, IL USA
[2]Committee on Genetics, Genomics & Systems Biology, University of Chicago, Chicago, IL USA
[3]Committee on Quantitative Methods in Social, Behavioral, and Health Sciences, University of Chicago, Chicago, IL USA
[4]Center for Health Statistics, Department of Medicine, University of Chicago, Chicago, IL USA
[5]Department of Neurology, University of Chicago, Chicago, IL USA
[6]Committee on Neurobiology, University of Chicago, Chicago, IL USA
[7]Division of Geriatric Medicine, Department of Medicine, Department of Community Health and Epidemiology, School of Health Administration, Halifax, NS Canada

★To whom correspondence should be addressed: e-mail: `ishanu@u chicago.edu`.

## CONTENTS

## LIST OF TABLES

## LIST OF FIGURES

### LIST OF ALGORITHMS

**a.** ZCoR AUC over time (95% confidence)

SI-Fig. 1: **Degradation of predictive performance with patient age.** With increasing patient age it become more difficult to distinguish age related cognitive decline from ADRD. This is reflected in the decreasing AUC with age, suggesting that comorbidity footprints associated with ADRD has confounding overlaps with conditions that arise more frequently as patients get older.

SI-Fig. 2: **Schematic of prediction pipeline.** Panels A, B and C show the sequential training steps, namely th e hyper-training of the PFSA models and the p-score dictionary, the pre-aggregation training of the four LGBM models, and the aggregation training of the final LGBM model respectively. Panel D shows the configuration of the trained pipeline in operation. In teh training steps, the filled box represents the component being trained in that step.

SI-Table I: High risk cohort definition based on known comorbidities of AD and related dementia

| ICD code | description |
|---|---|
| 250.0 | DMII wo cmp nt st uncntr |
| 250.02 | DMII wo cmp uncntrld |
| 252.0 | Hyperparathyroidism NOS |
| 252.02 | Sec hyprprthyrd nonrenal |
| 258.02 | Mult endo neop type IIA |
| 272.2 | Mixed hyperlipidemia |
| 278.0 | Obesity NOS |
| 278.01 | Morbid obesity |
| 278.02 | Overweight |
| 278.03 | Obesity hypovent synd |
| 296.2 | Depress psychosis-unspec |
| 296.21 | Depress psychosis-mild |
| 296.22 | Depressive psychosis-mod |
| 296.23 | Depress psychosis-severe |
| 296.24 | Depr psychos-sev w psych |
| 296.25 | Depr psychos-part remiss |
| 296.26 | Depr psychos-full remiss |
| 296.3 | Recurr depr psychos-unsp |
| 296.31 | Recurr depr psychos-mild |
| 296.32 | Recurr depr psychos-mod |
| 296.33 | Recur depr psych-severe |
| 296.34 | Rec depr psych-psychotic |
| 296.35 | Recur depr psyc-part rem |
| 296.36 | Recur depr psyc-full rem |
| 303.0 | Ac alcohol intox-unspec |
| 303.01 | Ac alcohol intox-contin |
| 303.02 | Ac alcohol intox-episod |
| 303.03 | Ac alcohol intox-remiss |
| 303.9 | Alcoh dep NEC/NOS-unspec |
| 303.91 | Alcoh dep NEC/NOS-contin |
| 303.92 | Alcoh dep NEC/NOS-episod |
| 303.93 | Alcoh dep NEC/NOS-remiss |
| 305.0 | Alcohol abuse-unspec |
| 305.01 | Alcohol abuse-continuous |
| 305.02 | Alcohol abuse-episodic |
| 305.03 | Alcohol abuse-in remiss |
| 401.0 | Malignant hypertension |
| 401.1 | Benign hypertension |
| 401.9 | Hypertension NOS |
| 402.0 | Mal hyp ht dis w/o hf |
| 402.01 | Mal hypert hrt dis w hf |
| 402.1 | Benign hyp ht dis w/o hf |
| 402.11 | Benign hyp ht dis w hf |
| 402.9 | Hyp hrt dis NOS w/o hf |
| 402.91 | Hyp ht dis NOS w ht fail |
| 403 | Hypertensive chronic kidney disease |
| 404 | Hypertensive heart and chronic kidney disease |
| 405.01 | Mal renovasc hypertens |
| 405.09 | Mal second hyperten NEC |
| 405.11 | Benign renovasc hyperten |
| 405.19 | Benign second hypert NEC |
| 405.91 | Renovasc hypertension |
| 405.99 | Second hypertension NEC |
| 427.31 | Atrial fibrillation |
| 440 | Atherosclerosis |
| E11 | Overweight |
| E66 | Type 2 diabetes mellitus |
| E78 | Disorders of lipoprotein metabolism lipidemias |
| F10.10 | Alcohol abuse uncomplicated |
| F10.159 | Alcohol abuse with alcohol-induced psychotic disorder unspecified |
| F10.20 | Alcohol dependence uncomplicated |
| F10.21 | Alcohol dependence in remission |
| F10.229 | Alcohol dependence with intoxication unspecified |
| F10.231 | Alcohol dependence with withdrawal delirium |
| F10.239 | Alcohol dependence with withdrawal unspecified |
| F10.27 | Alcohol dependence with alcohol-induced persisting dementia |
| F32.0 | Major depressive disorder single episode mild |
| F32.1 | Major depressive disorder single episode moderate |
| F32.2 | Major depressive disorder single episode severe without psychotic features |
| F32.3 | Major depressive disorder single episode severe with psychotic features |
| F32.4 | Major depressive disorder single episode in partial remission |

**SI-Table I – continued from previous page**

| ICD code | description |
|---|---|
| F32.5 | Major depressive disorder single episode in full remission |
| F32.8 | Premenstrual dysphoric disorder |
| F32.9 | Major depressive disorder single episode unspecified |
| F33.0 | Major depressive disorder recurrent mild |
| F33.1 | Major depressive disorder recurrent moderate |
| F33.2 | Major depressive disorder recurrent severe without psychotic features |
| F33.3 | Major depressive disorder recurrent severe with psychotic symptoms |
| F33.41 | Major depressive disorder recurrent in partial remission |
| F33.42 | Major depressive disorder recurrent in full remission |
| F33.9 | Major depressive disorder recurrent unspecified |
| G43.401 | Hemiplegic migraine not intractable with status migrainosus |
| I10 | Essential (primary) hypertension |
| I11.0 | Hypertensive heart disease with heart failure |
| I11.9 | Hypertensive heart disease without heart failure |
| I12.0 | Hypertensive chronic kidney disease with stage 5 chronic kidney disease or end stage renal disease |
| I12.9 | Hypertensive chronic kidney disease with stage 1 through stage 4 chronic kidney disease or unspecified chronic kidney disease |
| I13.0 | Hypertensive heart and chronic kidney disease with heart failure and stage 1 through stage 4 chronic kidney disease or unspecified chronic kidney disease |
| I13.10 | Hypertensive heart and chronic kidney disease without heart failure with stage 1 through stage 4 chronic kidney disease or unspecified chronic kidney disease |
| I13.11 | Hypertensive heart and chronic kidney disease without heart failure with stage 5 chronic kidney disease or end stage renal disease |
| I13.2 | Hypertensive heart and chronic kidney disease with heart failure and with stage 5 chronic kidney disease or end stage renal disease |
| I15.0 | Renovascular hypertension |
| I15.8 | Other secondary hypertension |
| I70.0 | Atherosclerosis of aorta |
| I70.1 | Atherosclerosis of renal artery |
| I70.209 | Unspecified atherosclerosis of native arteries of extremities unspecified extremity |
| I70.219 | Atherosclerosis of native arteries of extremities with intermittent claudication unspecified extremity |
| I70.229 | Atherosclerosis of native arteries of extremities with rest pain unspecified extremity |
| I70.25 | Atherosclerosis of native arteries of other extremities with ulceration |
| I70.269 | Atherosclerosis of native arteries of extremities with gangrene unspecified extremity |
| I70.299 | Other atherosclerosis of native arteries of extremities unspecified extremity |
| I70.399 | Other atherosclerosis of unspecified type of bypass graft(s) of the extremities unspecified extremity |
| I70.499 | Other atherosclerosis of autologous vein bypass graft(s) of the extremities unspecified extremity |
| I70.599 | Other atherosclerosis of nonautologous biological bypass graft(s) of the extremities unspecified extremity |
| I70.8 | Atherosclerosis of other arteries |
| I70.90 | Unspecified atherosclerosis |
| I70.92 | Chronic total occlusion of artery of the extremities |
| J12.0 | Adenoviral pneumonia |
| J12.1 | Respiratory syncytial virus pneumonia |
| J12.2 | Parainfluenza virus pneumonia |
| J12.81 | Pneumonia due to SARS-associated coronavirus |
| J12.89 | Other viral pneumonia |
| J12.9 | Viral pneumonia unspecified |
| J13 | Pneumonia due to Streptococcus pneumoniae |
| J14 | Pneumonia due to Hemophilus influenzae |
| J15.0 | Pneumonia due to Klebsiella pneumoniae |
| J15.1 | Pneumonia due to Pseudomonas |
| J15.20 | Pneumonia due to staphylococcus unspecified |
| J15.211 | Pneumonia due to Methicillin susceptible Staphylococcus aureus |
| J15.212 | Pneumonia due to Methicillin resistant Staphylococcus aureus |
| J15.29 | Pneumonia due to other staphylococcus |
| J15.3 | Pneumonia due to streptococcus group B |
| J15.4 | Pneumonia due to other streptococci |
| J15.5 | Pneumonia due to Escherichia coli |
| J15.6 | Pneumonia due to other Gram-negative bacteria |
| J15.7 | Pneumonia due to Mycoplasma pneumoniae |
| J15.8 | Pneumonia due to other specified bacteria |
| J15.9 | Unspecified bacterial pneumonia |
| J16.0 | Chlamydial pneumonia |
| J16.8 | Pneumonia due to other specified infectious organisms |
| J17 | Pneumonia in diseases classified elsewhere |
| K08.401 | Partial loss of teeth unspecified cause class I |
| K08.402 | Partial loss of teeth unspecified cause class II |
| S02.401A | Maxillary fracture unspecified side initial encounter for closed fracture |
| S02.401B | Maxillary fracture unspecified side initial encounter for open fracture |
| Y35.303A | Legal intervention involving unspecified blunt objects suspect injured initial encounter |

SI-Table II: Disease Categories With Detailed Set of ICD Codes Used

| Description | Constituent ICD9 Codes |
|---|---|
| Abnormal-Findings | R89.5 R92.2 R71.0 794.7 R87.620 I20.8 R97.8 R82.5 794.01 R93.9 R87.612 I25.3 I21.11 R85.616 R94.30 R87.611 I25.42 R94.39 R92.8 R88.0 R93.1 R82.2 R85.81 R94.09 I25.10 R82.4 794.6 R87.811 794.31 R87.810 I25.811 R79.1 R80.2 794.4 794.09 I25.2 R87.616 R73.02 R74.8 R87.9 794.14 R83.9 R75 I21.29 794.00 794.10 794.19 R85.612 R93.5 R87.628 I21.4 794.13 R94.112 R78.89 R94.120 R79.81 794.39 R85.615 794.16 R94.31 R73.09 R94.8 R87.624 794.17 R87.623 R89.9 I25.83 R85.82 R93.2 R94.4 R82.99 R87.619 R94.110 R85.613 794.5 R76.8 R80.3 794.2 R86.9 R91.8 I25.9 R94.01 R71.8 R82.3 R97.0 I24.1 I25.41 R87.621 794.11 R87.613 R93.3 R94.113 R94.121 R94.131 R94.5 I24.0 R85.619 R93.4 R97.2 R76.12 R85.614 794.02 R94.111 794.15 R70.0 I21.19 I20.0 R82.0 794.8 I20.1 R79.82 R94.130 R87.820 R87.625 R85.611 R87.622 R78.0 I24.8 794.30 R94.7 R93.0 R89.7 R87.614 I25.89 R74.0 R87.615 R85.9 794.9 I25.810 R90.81 R82.1 R93.7 I21.3 R73.01 R94.2 R97.1 I25.82 R87.610 R94.118 R92.0 R78.81 R91.1 R85.610 R81 I21.09 I25.812 R89.8 R76.11 R93.8 R94.6 794.12 |
| Acute-Bronchitis | 466.19 466.11 J20.9 I24.8 I25.10 I25.89 I20.8 I25.9 I25.811 I24.1 I25.2 I25.41 I25.810 I25.3 I21.11 466.0 I21.29 I21.3 I24.0 I25.82 I21.4 I25.42 I21.19 I20.0 I20.1 I21.09 I25.812 I25.83 |
| Allergic | 477.2 493.81 T50.995A J67.2 495.6 T78.03x 372.14 J67 J67.0 M13.89 J30.1 995.63 995.65 558.3 T45.0X1A M13.859 716.27 D29.30 D29.1 L27.2 477.9 495.5 493.22 D69.2 T78.00x 287.33 995.60 J45.31 J45.51 D29.20 J67.7 T78.09xA D29.22 M13.80 J30.9 T78.08x 287.8 H10.45 B44.81 716.20 995.61 T78.05xA 493.92 693.1 493.90 T78.40x J45.20 493.82 J45.40 D69.42 495.7 J67.5 493.20 D69.49 J45.32 287.32 708.0 H65.119 995.64 D69.1 J45.21 D69.6 M13.819 716.23 495.4 995.67 287.1 T78.08xA T78.00xA 477.0 493.02 525.66 T78.02xA J67.1 D69.3 T78.04x T78.2xxA D29.4 716.25 T78.07xA 716.26 T78.07x M13.88 J67.3 495.9 J45.30 493.21 477 495.2 995.62 T78.40xA 995.27 287.2 495.8 495 287.5 995.0 493 T78.05x L50.0 493.11 J45.902 D29.0 J45.990 287.9 J45 D29.21 J30.0 963.0 495.1 D29.32 L25.9 J44.9 J44.0 477.1 M13.879 493.01 J45.41 T50.995 J45.998 692.9 M13.849 995.66 D69.8 995.69 T78.04xA J30 495.3 M13.869 287.30 J45.991 J44.1 995.3 287.4 J45.52 287.0 381.06 716.21 J45.901 J67.4 287.39 493.91 373.32 287.31 T78.06xA J30.89 287 K08.55 K52.2 D29.31 J45.50 495.0 J67.6 D69.9 D29.8 T78.02x 716.24 477.8 381.05 D29 493.12 T78.03xA J67.9 716.22 T78.2xx J30.5 999.4 493.00 M13.829 T78.01x T78.06x 493.10 518.6 716.28 J30.2 H01.119 995.68 M13.839 D69.0 T78.09x 381.04 D29.9 T78.01xA 716.29 J30.81 J45.22 J45.42 T45.0X1 J45.909 D69.41 J67.8 |
| Cardiovascular | I35.0 I48.0 I25.728 444.8 P29.38 I94 I63.212 402.00 I70.669 440.30 I89.9 I60.9 I20.1 413.9 I24.1 I80.3 415.1 I77.811 785.9 I69.319 I69.339 I82.5Z9 R04.1 429.6 G43.619 I82.503 I82.611 I70.512 I75.011 I69.834 I70.628 K64.9 I89.0 I21.09 428.42 447.5 442.8 I70.792 454.0 I70.318 I50.83 I70.744 405.0 426.2 455 I70.442 455.3 I82.B29 I12 415 433.8 I69.320 I27.81 444.21 I70.735 I82.602 I67.89 441.4 425.4 I35.9 I70.693 I69.234 I65.23 427.2 I70.244 I49.02 I82.91 P29.89 I70.719 I69.131 I36.8 I60.8 I60.11 442.82 I69.852 I75.029 438.22 I69.120 I70.641 I60.2 426.51 I70.302 417.9 I63.012 R04.2 R00.2 427.31 I25.718 I05.8 I70.791 I89 426.50 I63.349 I49.49 444.89 I63.213 I83.12 I77.0 I82.433 I70.608 P29 I47.1 428 I70.348 I82.C29 I82.532 I69.322 I63.311 I65.03 I82.291 I70.498 I97.791 I69.859 I70.644 I82.441 I63.413 I70.362 405.11 I62 I87.012 I80.292 411.1 433 I70.532 I97.638 I87.091 I69.998 I07.8 I11.9 I69.390 445.01 I69.223 I37.2 I87.319 I69.932 I70.208 I82.492 I82.891 I63.233 I70.319 I70.65 I70.341 435 441.9 I40.0 I63.539 K64.4 I95.89 I69.028 I82.5Y9 I50.32 I70.731 I70.768 458.0 G43.601 I21.01 410.71 429.3 I69.398 I87.399 441.0 I70.421 I73 I70.729 E86.0 I69.231 I28.8 I13.0 I70.568 I42.2 I63.431 411.89 I70.709 438.21 438.53 I70 I50.84 I42.0 I70.212 I77.74 I08.0 P29.2 455.0 410.1 416.1 I24.8 433.10 I70.393 413.1 I70.561 I87.092 I83.218 427.60 453.82 I70.534 I97.130 I97.821 I97.711 I63.543 I02.0 405.01 I69.364 I07 I83.92 I69.928 I69.214 I70.418 346.63 I06.8 410.60 I31.2 I70.433 I75.013 I70.219 I70.431 I75.023 I82.403 I83.10 I25.730 415.3 I13.2 I69.814 I50.812 I79 429.89 I69.833 I23 410.20 I77.6 I69.165 I60.52 459.81 I47 410.10 I83.201 I82.419 I69.033 I15.9 404.0 I69.213 441.7 I82.412 I69.261 I82.432 417 I24.9 I69.315 438.19 I72.8 I36.2 I97.648 I70.539 I56 I69.392 I63.512 I69.820 I77.70 I63.00 I70.439 I70.643 433.80 276.50 459.3 I70.522 I69.121 433.21 426.4 346.61 I61 I75.021 I88.8 I55 785.1 I01.9 I25.729 I63.20 454.9 I63.521 I70.769 I69.334 I72.4 I50.31 438.50 434 I70.349 I48.2 435.1 I80.219 I10 I69.115 I70.544 414.00 I69.010 I66.12 445.8 I30.9 I27.1 I80.293 432 I97.190 I69.864 I86 I34.1 410.8 438.9 I63.133 I25.89 I63.032 I34.8 I77.4 I82.4Z2 I87.391 411.0 I87.312 440.4 I63.232 I87.009 R00.8 I69.012 I70.468 Q82.5 I21.3 I87.303 I25.6 I69.842 290.41 455.8 I69.042 I82.C23 I87.099 I69.314 I74.10 I77.812 I69.393 I42.5 I69.141 I70.491 I82.210 785.50 I69.351 I82.513 785 I69.252 785.52 I77.819 415.11 I80.13 I71.5 I82.422 I42.6 P29.30 I69.031 I70.222 I31.4 I82.603 I82.729 437 I70.621 I70.509 I70.763 I95 I38 433.3 I69.133 K55.0 426.6 I80.212 I70.723 I70.229 404.10 I83.029 437.3 426.13 I65.22 I75.012 I80.202 I48.3 I27.83 I69.815 I75.022 I27.20 I70.711 453 I65.9 I80.233 410.01 I70.291 440.2 I63.29 404.11 I82.A21 R04.89 I08.3 785.0 438.10 I70.469 I91 I69.163 447.2 I50.1 414.4 I83.002 410.5 I70.243 I63.039 I70.249 I70.269 I80.211 I86.8 I70.499 I69.363 404.02 I22.0 I66.01 458.2 I83.213 I83.028 I74.2 428.0 I70.545 442.83 426.3 I69.362 I31 I70.398 I81 I69.810 I70.261 443.22 I70.1 I70.335 I92 I77.89 I60.00 I83.203 I72.3 I63.532 410.2 I71.02 I70.331 441.03 I69.262 410.6 I65.29 I54 453.71 I69.293 438.81 I69.110 I15 R65.21 428.31 453.89 I69.051 I67.7 I69.062 I69.111 I69.349 454.1 453.1 I63.39 I82.1 I66.03 I72 I46.8 426.12 I45.81 I73.1 I70.634 I69.918 I95.1 I87.393 410.02 410.0 I70.332 I82.549 443.29 I77.5 440 I31.8 I21 I82.B23 I11.0 433.01 455.7 403.00 I66.02 438.14 I06.0 G45.0 I97.51 I69.112 I70.369 I82.5Y2 I45.2 I74.01 415.19 I87.311 I63.22 I84 I20 I70.203 I83.009 I87 I71 433.1 I06.1 456.4 I69.822 I99 I33.0 I70.218 438.5 I60.10 I69.254 413 I51.2 I69.065 I59 410.62 I07.2 I80.8 I96 453.76 437.2 435.3 414.11 I83.019 I09.2 I83.93 426.81 I63.323 410.72 435.2 I63.441 I83.011 410.31 I82.B11 410.81 I27.23 438.41 I75.81 404.9 I13.11 I66.09 I95.0 I83.022 438.4 I82.431 I70.202 I70.342 I70.648 I23.3 I69.265 437.9 I82.812 I82.439 I63.59 I97.120 I69.249 I07.9 455.6 I82.501 I82.A23 I74.09 I63.329 I70.722 I82.423 I87.9 I82.C21 414.02 I42.9 I69.093 404 I24 455.9 I32 I70.798 I44.60 441.02 R04 I50.40 I87.302 I82.509 441 I77.3 I97.630 I69.993 I53 I83 I82.499 I69.034 I70.402 I60.01 458.1 438.82 I24.0 458 I69.144 R01.2 I70.231 I65.1 I70.429 I70.612 I50.33 I82.612 I63.422 I57 433.2 I69.221 I95.9 I78 I97.111 I04 I69.963 779.82 410.52 410.3 I70.298 438.20 I48 I50.82 I70.698 I70.533 I69.915 I69.052 425.9 I69.243 I80.229 427.32 438 410.41 I69.092 413.0 290.4 I27.22 I69.043 442.83 453.86 557 453.75 445.0 I23.2 438.1 458.21 I70.303 I70.541 I70.201 I69.154 411 429.1 I51.9 I60.7 I69.244 427.41 I72.6 I16.0 290.42 I70.793 I82.491 I69.022 443.8 I51.5 I26.09 I83.225 I37.9 I66.13 I70.329 I69.191 I70.638 444.2 I36.1 I69.91 R01.1 I08.1 I70.333 I70.462 I82.401 I82.A19 I63.519 K55.9 I82.4Y3 I82.5Z3 I33 410.7 I45.4 I69.333 I83.219 I69.198 I70.399 426.54 I25.710 442.81 557.9 I80.231 I82.443 I80.01 P29.12 447.3 I70.245 410.70 433.91 I69.812 I50.23 448.9 434.9 I26.99 427.5 I77.77 453.6 I69.118 I35 441.3 I82.C22 I08.8 414.07 I28.0 438.89 I97.42 I65.09 I37.0 I67.848 I70.449 I77.2 411.8 I70.603 G45.4 453.81 I44.5 426.0 I06.2 I74.5 427.9 427.42 I70.213 I61.2 444.9 I62.1 F01.51 I70.701 456.8 456 I87.013 I01.8 440.3 453.83 I69.242 I42.1 I83.811 I69.069 I69.313 I83.023 I45.9 I69.013 I69.813 453.77 I21.11 I25.82 I87.339 I70.749 I83.892 I23.5 I97.611 I66.3 I63.449 453.72 I88.0 405.9 I69.059 I40 I69.943 I70.748 I21.02 428.32 447.70 I78.8 403.9 I82.521 I69.321 405.09 453.40 I63.9 I69.854 453.42 I40.8 I69.162 426.1 I69.192 I83.222 I87.033 I25.83 I74.3 I82.523 458.9 K64.8 438.83 I42.7 I70.613 I83.221 I82.609 I82.703 I70.07 I83.011 I69.942 437.6 405.1 I60.6 I69.342 I70.728 405.99 I69.831 I69.30 I12.9 438.3 447.71 441.00 I69.828 434.91 I09.89 I23.4 I83.024 I47.0 I82.543 I50.810 I70.293 I82.211 405.19 I86.2 I69.954 I70.409 I83.812 I97.622 437.0 I25.799 403.01 I75.89 I50.20 I63.533 I78.1 428.33 I82.A22 I79.0 414 402 434.01 447.89 I69.964 I70.642 I63.542 I47.2 I87.309 I63.339 I77.72 I70.363 414.06 404.92 I71.2 441.1 I82.522 I20.8 402.11 I82.819 I65.8 I36.9 I99.9 785.3 429 I69.344 I97.811 I63.113 443.21 I69.020 I70.328 I72.9 I62.01 I70.613 447.72 I82.542 785.51 I69.132 I83.015 I82.493 I67.4 426.52 425.7 I80.00 I70.338 I63.30 I78.9 I74.19 I87.392 442.9 I82.449 I69.862 403.90 402.10 I25.758 402.1 I27.21 I83.202 I86.0 I70.401 I87.093 I49.8 I95.3 433.11 I02.9 425.2 I70.519 412 I70.739 I70.702 404.03 429.5 I88.9 410 I82.B21 410.22 I69.169 I70.368 429.82 I82.619 I82.592 I60 437.8 I80.10 414.0 433.9 I72.5 I70.591 404.12 I69.219 I73.9 I46 I63.50 434.1 I82.599 I74.9 I83.018 I82.4Z9 I70.392 I66.11 I63.531 I70.435 I22.2 I65.02 I82.541 425.0 I69.298 I63.40 I69.263 I95.81 453.8 I82.409 R03.1 I25.759 I87.1 I69.353 434.01 557.1 I25.2 I70.598 I70.601 I79.1 R57.8 I63.131 I22 414.03 I46.2 I19 I63.439 403.91 I69.064 429.8 426.9 I63.09 I80.03 I45.89 I97 I82.531 429.79 I13 I69.032 I22.8 I26.90 434.10 415.12 I70.639 I69.232 I51.0 I60.51 I63.411 I97.131 I69.891 I70.8 I44.0 I39 |
| Cataract | 366.34 366.14 366.19 366.02 366.21 366.31 366.10 366.23 366.32 I24.8 I25.10 366.16 366.18 I25.89 I20.8 366.15 366.52 366.00 366.51 I25.9 366.13 I25.811 366.03 I24.1 I25.2 366.17 I25.41 I25.810 I25.3 I21.11 366.04 366.12 366.11 366.9 366.30 366.33 366.01 366.43 I21.29 366.09 366.53 I21.3 366.44 I24.0 I25.82 I21.4 I25.42 366.8 366.20 366.50 366.22 I21.19 I20.0 366.45 366.46 I20.1 I21.09 I25.812 I25.83 366.42 |

| | |
|---|---|
| **CNS** | 346.32 335.10 G40.019 G47.54 G43.609 G44.029 G16 344 349.2 G47.20 G45.0 G21.4 347.00 G81.11 346.82 G95.19 G47.27 345.70 347.11 G25.89 341.21 322 344.30 346.93 G04.2 G40.909 G40.219 G43.619 G80.0 327.15 G43.011 333.1 349.82 327.51 G83.10 G12 G47.35 327.32 G03.8 327.41 337.00 G46.2 G96.12 G05.4 G43.101 G24.2 G47.29 G12.21 345.00 G93.1 346.72 G13.1 G95.9 331.7 G31.01 343.8 G40.209 346.40 334.9 G30.1 G44.309 G83.13 321.0 344.89 G47.23 G44.009 G11.9 324.9 G43.D0 G81.90 346.20 G03.1 333.84 320.1 348.30 G30 G81.93 G95.0 G37.5 G47.12 G25.71 327.29 G40.201 346.41 G06.2 G37.9 767.0 G40.813 P91.63 G93.49 344.9 345.01 334 G47.52 322.9 G44.41 336.1 G44.82 335.11 348.39 G47.22 G37.2 G46.0 327.11 G04.89 G22 331 G41 G43.509 G83.30 327.30 G93.7 333.2 333.5 G47.419 344.41 P91.2 P91 328 349.1 G90.529 G25.79 344.5 G40.111 337.20 G83.89 G37.1 348.4 G23.2 G25.1 G44.81 G00.2 G40.804 G40.B09 336 G44.83 G47.00 G82.50 344.60 337 343 G47.63 327.44 G44.40 336.9 G13.8 P91.1 333.82 G47.59 G31.1 G82.20 334.2 333.71 G35 G21.0 G24.1 G25.3 G23 G46.5 322.0 327.8 779.2 333.85 G25.0 G81.01 344.04 G43.C1 G44.53 G40.814 G00.3 G43.111 342.02 333.89 G23.1 G90.59 G80.2 F51.05 G91.0 G40.823 343.4 344.01 335.20 G40.411 G44.221 G47.14 G12.25 G04 343.3 G43.A0 333.79 327.22 346.01 G18 G03 G83.9 G40.409 G44.001 331.82 323.42 323.72 G82.53 337.3 G93.40 G99.0 G43.719 E75.4 G40.009 G47.30 G47.10 G40.802 336.0 327.01 G40.B01 P91.60 323.63 340 G43.B1 345.91 P52.8 G12.8 G43.419 G07 G93.9 G40.001 327.40 G44.209 337.1 G25.4 G99.2 333 333.72 323.81 327.23 323.02 344.2 G12.24 333.4 327.59 G12.23 G43.601 331.2 333.3 G83.4 G44.219 G40.101 B45.1 330.0 G44 G83.23 G93.2 G43.001 G43.C0 327.49 G21.8 G46.7 G30.9 G40.419 G82.52 G40.A01 345.61 347 G44.321 G11.8 G00.9 330.1 G43.B0 G44.011 323.1 G31.81 G43.401 320.89 G83.20 346.80 G47.9 G43.911 G47.62 348.9 G13.2 346.90 G47.31 G90.09 335 342.82 327.20 342.11 G11.2 327.19 346.02 341 G40.119 G47.09 322.1 G04.02 G19 348.5 G43.829 323.82 333.99 342.12 345.3 G40.011 G20 G45.9 G40.89 P91.0 346.42 346.60 335.23 G03.9 G40.309 336.3 333.90 349 G04.01 346.00 327.14 342 G05.3 G44.059 321.2 G31.89 G43.701 G37 G43.511 G40.822 G43.919 G44.51 G47.36 345.41 G31.2 G13.0 G91.1 344.31 G00.0 320.2 G25.81 344.61 G44.329 G47.429 G40.824 G80.8 346.63 G46 G30.8 G43.831 G80.3 348.2 333.83 G45.4 G13 320.0 345.10 344.40 341.1 G37.3 G31 G14 333.94 G24.3 G36.9 G45.2 349.39 G24.8 G47 G04.31 G05 G11.4 341.20 332.1 344.1 G47.39 G83.5 G46.1 G25.70 G21.19 G00.1 343.9 G25.61 G40.311 G21 331.11 G32.89 G47.51 332.0 335.22 G44.091 345.80 331.81 G36.8 G00.8 G46.8 G47.34 P91.88 G43.D1 323.71 330.9 G04.90 331.89 G06.1 G46.4 348.1 346.30 G43.019 G82.54 337.29 335.8 G04.1 337.22 G43.839 344.02 327.43 G12.1 G43.901 G47.69 G10 G47.421 G47.50 G44.84 G31.83 346.62 G47.01 345 G47.19 G44.031 G25 320.7 323.51 G33 G06.0 G32.81 346.50 G83.11 G27 337.09 G24.09 322.2 344.00 346.61 G43.801 327.37 G31.85 327.42 G36 346.31 G46.3 G24.9 G40.211 G24.01 327.26 G44.049 341.9 G44.039 327.33 346.83 334.1 335.29 344.09 323.2 346.13 327.12 G47.37 327.36 327.02 G43.909 G25.5 G40.803 346.12 G12.20 327.31 G37.8 345.2 327.10 G44.311 G43.409 344.42 343.0 G06 323.01 333.93 346.81 G40.509 325 G44.211 G43 G81.00 G93.0 341.8 335.0 G09 G96.9 320.81 G44.229 G04.32 G24.4 346.11 G25.69 333.6 G37.4 G04.81 G44.201 345.81 332 G12.22 346.51 346.33 331.5 342.90 G44.099 G40.401 G95.89 346.10 F51.13 G43.501 G01 334.4 324 327.13 G44.051 345.71 327.53 320.3 345.60 G12.0 G04.91 G44.85 G43.821 G97.41 G11.0 346.21 P91.62 G39 345.51 321.1 G96.11 323.62 G80.1 G43.009 G31.82 344.81 327 G43.519 346.23 G43.119 336.8 345.40 345.50 333.0 G81.10 348.89 326 345.90 327.24 333.92 G29 334.0 G40.301 G47.24 G43.711 G44.021 G46.6 G40.319 333.91 G04.30 G43.A1 G31.84 G43.411 G32.0 337.01 G04.00 G21.9 P91.61 G47.26 346.43 G44.52 G83.21 G38 344.32 327.34 346.71 P91.811 G11.1 G40.A19 G40.B19 324.1 G40.109 G43.109 G90.4 G15 G47.33 342.92 G40.901 G40.811 G93.89 G93.6 E75.23 G21.3 G94 342.01 G45 G37.0 G12.29 G47.61 334.8 343.2 349.9 342.80 G93.41 G25.83 G43.709 G43.819 335.9 G43.811 346.53 347.01 G47.8 G97.82 G34 342.81 342.91 G44.10 G43.611 G03.2 G31.9 G44.1 327.21 323.0 G31.11 G32.89 G47.51 332.0 335.22 323.41 320.82 G40.B11 P91.9 337.9 G00 335.24 G47.411 321.4 348.31 346.91 327.27 330.2 348 327.00 348.81 G03.0 346 G44.041 321.3 G31.09 346.03 G28 G21.11 G25.2 G02 349.0 327.52 G40.A09 G40.821 349.89 G17 346.22 327.35 330.8 331.83 G81.13 G97.1 G23.9 329 G90.01 G40.911 333.81 G21.0' P91.5 G43.809 335.19 G45.1 336.2 G24.02 335.21 346.92 G23.8 P91.4 341.0 G83.0 327.25 G47.25 346.52 G21.2 349.31 G25.9 G40.812 G24 337.21 331.4 G47.53 346.70 G91.2 345.11 330 327.09 344.03 G96.8 323.9 G44.019 G47.13 323.52 G44.301 324.0 G90.9 G42 346.73 320.9 G44.59 G25.82 G40.919 330.3 327.39 321.8 331.0 331.9 G30.0 G36.1 779.1 G81.03 334.3 G08 G47.21 331.3 G40 G11.3 G32 G04.39 342.10 G92 G36.0 323 G82.51 G45.3 G12.9 349.81 P91.819 321 G47.11 G40.A11 348.0 G93.81 G11 G40.801 G26 G80.9 342.00 G40.501 320 P91.3 G44.89 G90.519 G93.5 G81.91 G24.5 G45.8 331.19 323.61 G96.0 343.1 |
| **Development** | 520.4 742.59 764.96 743.66 750.19 M26.36 756.11 749.01 748.1 524.33 740 740.1 752.35 741.0 743.53 524.79 748.5 752.42 589.0 Q18.1 Q26.2 Q93.4 764.98 743.43 744.21 Q63.8 753.23 750.4 P92 N13.70 746.82 747.40 743.57 743.41 742.5 752 743.34 749.14 524.35 744.82 Q05.2 M26.9 Q55.63 Q10.1 752.43 748.3 756.16 M26.89 743.31 593.71 M26.221 M26.72 313.23 Q31.1 743.49 F94.0 750.24 747.41 Q35.7 Q23.0 Q15.8 Q61.3 749.1 Q16.9 759.0 743.12 750.27 749.03 747.49 Q61.8 P92.4 743.65 750.21 Q20.0 Q87.81 Q62.4 758.5 764.03 524.73 Q17.1 756.13 524.31 Q34.8 Q18.4 Q27.8 520.8 743.64 Q21.2 741.9 751.61 Q76.5 524.27 Q05.7 Q61.01 746.87 Q10.0 Q25.4 Q36.9 764.06 752.36 743.11 764.91 Q28.3 747 756.14 743.39 759.89 746.6 753.12 750 750.10 Q62.31 M26.39 750.6 Q51.4 Q56.3 Q24.8 747.82 745.7 752.65 750.16 749.02 Q28.8 P05.14 743.35 Q20.9 Q64.9 P05.08 749.2 742.9 747.10 756.12 Q16.4 743.44 K00.4 P05.15 Q22.5 747.20 745.19 753.4 764.07 Q51.3 N28.83 M26.56 Q21.1 744.8 M26.220 747.8 N13.722 Q23.2 752.40 Q00.1 Q87.1 Q06.4 764.92 P05.06 524.56 743.46 Q55.22 745.11 744.0 Q51.2 K00.2 M26.52 Q76.49 752.89 Q52.12 746.0 Q99.2 E30.1 P92.2 M26.54 745.12 Q02 753.10 745.8 Q33.4 Q24.2 764.12 Q24.6 754.0 315.8 744.89 307.6 Q55.23 750.1 759.4 M26.73 Q00.2 Q40.8 752.64 M26.211 593.7 Q64.10 758.32 Q33.6 524.3 Q26.5 Q18.2 Q55.64 746.5 P05.9 764.00 744 752.34 752.9 Q43.4 524.57 Q38.2 759.1 747.2 315.9 Q13.4 759.3 524.22 Q16.2 315.4 743.63 752.51 Q27.9 745.3 524.32 743.54 Q20.1 F93.9 Q12.0 746.8 524.8 P05.18 748.61 P92.9 Q42.9 Q52.0 743.56 742.1 Q20.8 Q89.3 524.21 Q51.0 749.20 753.3 Q18.7 315 Q51.811 747.3 Q03.8 748.60 744.09 Q61.9 Q93.81 745.69 741.92 759.2 Q05.8 746.2 747.81 750.29 741.02 Q76.0 744.04 743.36 764.24 P92.8 P05.10 749.11 Q12.4 524.55 Q54.4 746.4 520.1 743.59 764.10 Q37.8 M26.33 313.2 P05.13 524.7 P92.01 745.61 752.11 P05.16 M26.34 Q24.3 751.5 524.74 Q06.2 746.00 752.10 Q50.6 Q21.9 752.49 745.4 758.31 752.33 K00.0 E30.0 753.15 307.7 758.1 746.1 Q04.8 764.13 M26.30 752.41 Q10.6 744.23 759.9 764.17 764.90 Q14.0 Q13.0 764.99 593.70 Q43.0 753.21 P05.07 744.05 745.6 Q98.4 744.47 Q52.9 750.3 743.45 Q33.9 764.97 753.7 741.93 764 Q23.1 Q30.0 747.1 752.6 Q15.0 593.73 741.91 M26.79 742.51 524.81 Q51.818 764.21 Q33.0 743.1 Q22.0 Q10.7 Q51.5 752.69 752.1 Q22.1 748.6 Q12.1 750.8 F93.8 Q36.0 745.0 751.8 752.4 Q34.9 747.83 747.21 754.1 746.84 753.5 520.9 524.30 747.60 747.61 R62.51 K00.9 743.52 742.2 Q12.9 753.0 Q51.6 743.9 743.5 759.7 Q89.7 744.03 743.03 753.22 Q38.0 Q38.3 759.81 M26.23 745.9 M26.24 744.02 524.82 P92.5 Q01.9 751.60 Q26.8 Q17.0 Q41.9 524.20 Q97.1 764.20 Q52.3 Q13.5 Q38.4 752.8 752.3 M26.4 746.3 M26.55 747.22 520 743.62 741.00 M26.57 749.04 744.9 Q17.8 756.2 758.4 Q45.1 M26.32 F98.0 P92.6 P05.12 Q17.2 524.26 524.50 748 P92.09 743.2 P05.05 741 756.4 Q10.3 N27.0 747.9 309.21 Q05.0 745 P92.3 Q14.1 742.3 Q62.39 752.47 Q51.810 Q17.3 740.2 Q14.2 Q11.1 589.1 746.09 Q45.9 743.10 747.89 K00.5 751.0 Q15.9 Q64.39 750.11 Q68.0 750.13 Q62.10 749.24 751.9 313.8 743.42 Q77.1 750.26 Q22.2 744.83 524.5 593.0 259.1 747.4 744.01 524.29 744.5 524.53 K00.6 Q04.3 Q92.8 Q12.3 752.44 752.0 751.69 Q89.01 764.08 743.6 764.15 F88 315.5 748.69 589.9 Q11.2 Q55.8 752.7 764.16 747.42 743.33 744.81 Q25.2 741.01 Q62.12 Q05.5 Q89.4 744.24 Q13.3 Q95.0 P05.02 N27.1 313.9 Q38.5 748.4 742.4 744.4 Q50.4 524.54 P05.9 P29.3 P05.01 524 Q99.9 749.12 758.9 758.33 764.11 753.8 Q07.9 745.5 748.0 747.63 589 758.6 P05.17 764.05 Q38.6 Q26.3 Q26.9 M26.70 524.34 747.0 Q20.5 M26.59 764.02 Q06.8 743.20 Q00.0 764.29 M26.213 Q33.1 752.5 Q76.419 Q91.3 K00.1 746.28 751.4 747.5 752.62 M26.25 524.70 748.8 Q43.3 743.58 Q35.9 749.10 Q21.0 752.2 N13.729 752.63 753.17 Q18.5 Q89.2 R62.50 M26.212 Q16.1 Q61.4 Q20.3 524.24 747.62 747.11 749.2 524.37 764.95 743.22 744.22 741.03 Q24.4 Q22.3 M26.82 744.84 Q16.0 Q24.5 753.1 259.0 748.2 752.46 747.64 764.25 764.01 Q39.5 524.72 756.17 743.21 756.15 750.23 Q51.820 Q38.1 Q18.8 Q60.2 P92.1 524.59 750.25 750.5 749.00 752.39 M26.50 752.45 744.3 743.32 F93.0 Q45.8 313.89 758.81 750.0 Q64.4 Q96.9 742 746.85 743.30 313 743.06 Q14.8 R62.7 743.3 Q05.4 751.6 746.89 758.39 749.25 R62.52 752.19 N13.721 752.52 744.1 Q76.2 764.04 Q16.3 745.2 P05.04 748.9 764.14 764.94 520.7 749.21 753.19 750.7 Q25.0 764.27 F98.1 752.31 Q20.4 P05.03 Q91.7 759.83 746.81 743.48 749.0 Q50.01 524.9 Q40.1 745.10 524.76 P05.11 524.89 753.2 750.2 Q27.31 R62.59 744.41 744.46 751.1 740.0 Q67.4 749 K00.8 F82 Q64.79 Q23.4 764.93 Q27.0 753.11 Q11.0 M26.74 Q18.0 Q55.62 M26.37 Q21.3 Q22.9 Q25.1 524.4 745.60 750.12 752.32 742.53 Q37.9 Q44.2 524.71 R62.0 Q51.828 N27.9 753.6 743.8 593.72 524.54 744.00 524.25 750.9 Q61.5 Q30.8 Q23.3 751.2 524.28 746.86 756.10 524.52 M26.35 746.83 743.61 520.2 764.19 Q28.9 744.42 Q18.9 746.01 749.22 744.49 M26.81 743.37 751.7 747.29 P05.2 520.3 524.36 750.22 764.22 Q64.0 Q76.1 742.0 747.6 744.2 Q14.3 746.7 Q53.9 749.13 743.55 Q62.11 Q80.0 515.6 Q16.5 Q27.2 Q05.1 Q52.11 752.81 758.2 520.6 M26.31 524.39 Q18.3 Q93.89 M26.20 Q38.7 755.55 524.75 764.26 M26.53 744.4 758.7 756 524.23 R62 K00.7 Q40.0 Q05.6 524.2 Q13.81 K00.3 745.1 756.19 752.61 Q31.0 758 520.0 741.90 746.9 Q93.88 Q89.9 753.29 Q13.1 Q18.6 F81.9 746.02 743.00 Q44.1 743.47 747.69 Q39.1 Q38.8 753 764.23 764.18 753.20 Q27.32 764.09 Q17.9 Q23.8 Q89.1 Q13.89 M26.71 750.15 753.9 743.69 744.29 P05.00 Q40.2 |

| | |
|---|---|
| Digestive | 533.60 K63.81 531.30 K13 K42.1 K52.81 K11.4 K68.11 R15.2 K94.13 K42.9 K52.3 R10.811 R10.823 K26.2 578 K22.719 K64.9 532 K29.81 K76.2 K57.53 K59.02 532.50 K80.42 K57.80 K80.34 K08.123 534 K50.912 K56.601 K58 579.9 531.91 K76 K55.022 K55.30 K56.600 K08.133 K91.840 537.82 534.5 K08.51 K45 527 K80.30 K67 K63.3 K22.2 538 787.5 K26.5 K02 K66 564 K80.13 K22.6 R13.14 K51.211 K39 K05.00 K70.31 K00.4 K08.403 K13.29 K43.7 K70 541 R10.816 K09.0 K59.31 K85.82 532.11 K55.029 K00.2 K60.2 K92 K80.20 K08.122 579.1 532.91 K03.9 K76.5 R19.31 533.61 577.1 K50.919 K25.9 533.30 K35.2 534.2 K56.1 K63.89 K37 K57.40 533.91 K51.919 K05.212 K92.81 K13.79 R19.36 K62.1 K52.831 K51.40 531.0 K51.312 K50.812 K91.2 532.51 787.4 K40.40 K04.01 P77.3 K90.0 K22.3 K84 532.9 K64.4 K50.911 K69 K74.2 K92.9 K03 K64.1 R19.06 K08.119 K27.9 K57.30 K01 K08.414 K08.530 K21 K43.1 534.60 K08.0 K05.30 K80.44 562.12 K27.3 K56.699 K11.23 K50.118 R19.11 531.2 K13.1 K08.9 534.00 530.5 K11.9 527.8 K40.30 569.44 K62.3 K22.70 K08.433 K55.042 K70.30 K64.0 K94.01 K65 R10.2 533.31 K04.99 577.8 K94.23 K08.22 K41.41 K00.9 568.89 K41 K30 K08.52 K76.4 E16.9 K80.19 K02.7 K70.0 578.1 K27.6 K25.7 J86.0 K86 K94.39 K10 K05.323 P78.82 569.41 K32 K29.71 K04.3 K91.0 K58.8 K59.04 787 K57.12 K08.23 K95.89 K60.3 787.7 K26.9 K28.2 K14.1 K00 533.51 K08.494 K57.20 456.2 K04.90 K52.29 K56.0 K28.0 K13.23 530.82 K42 531.41 K29 R19.30 K27 K80.31 K59.01 777.2 R14.2 R14 R19.07 R10.824 534.9 K51.80 K22.10 R11.12 R18.8 K55.069 K28.4 787.21 K31.0 K55.21 K51.913 K38.8 R19.7 K46.0 K55.8 K80.51 K63.1 K57.31 536.1 K85.80 K85.30 K56.50 537.84 K94.29 532.31 534.30 K52.22 K40 R10.84 K73.2 K81.1 K80.66 K71.4 530.8 K40.10 K59.8 K71 K58.9 K80.00 K50.814 527.2 787.23 K54 K08.431 532.40 569.85 531.31 564.01 537.83 569.87 R11.0 K14.3 537 K62.5 564.09 R14.3 K51.314 K22 K86.1 K91.5 K05.211 K71.50 K08.104 530.9 534.51 K74.1 534.70 K80.11 K62.89 K76.3 K89 R13 K19 K62.6 569.2 K31.1 K65.2 K91.31 532.71 K11.5 K60.4 K57.21 K08.422 K71.9 K51.514 K59.39 K41.10 K08.531 K03.81 K75.3 K51.411 530.11 531.01 K22.9 K72 K64 K08.423 K28.5 K29.01 K08.21 534.40 K31.89 530.2 K27.5 K51.512 K29.41 K22.4 K27.1 K02.61 R10.821 K51.518 K85.10 537.89 K73 K27.2 K38 558 536 536.8 541.0 K94.03 K86.0 K57.13 533.4 K28.9 K38.1 K93 K51.011 R10.826 K50.818 456.1 K51.912 K08.20 568.82 K61.0 579.2 K14.9 K82.1 K06.8 K08.50 R19.6 K51.218 565.0 K52.89 R13.19 K51.911 568.81 P76.1 K57.11 558.41 K08.424 K55.021 R10.812 K91.72 K50.013 K70.9 K51.219 K59.09 K94.21 K05.312 531.60 K09.1 777.6 R15.0 K29.20 562.13 K91.83 K56.691 K05.11 K31.83 K60.0 K08.412 K55.032 K13.70 536.0 K08.101 K51.812 568.9 K44.0 K34 K55.039 R16 564.9 K65.0 K12 K04.02 K08.55 P78.1 K08.192 532.30 543 K71.0 K65.3 K08.59 K66.9 K35.89 560.3 K80.70 577.2 K31.84 K08.499 K41.21 P78.83 K14.6 K08.111 K41.00 K65.1 K86.8 568.0 K26 K51.311 K38.0 K76.81 787.91 K29.51 K07 K94.11 K38.9 578.9 K03.89 K55.049 543.9 K85.31 K52.82 R19.02 K08.82 566 R10.30 565.1 K14.0 K85.22 531.90 K29.90 527.4 K94.10 K80.80 527.6 K08.409 K72.91 K49 R10.13 R19.05 532.6 K12.2 K41.01 K60.1 K29.00 K55.9 R15.9 K66.0 K08.124 K08.421 K02.63 K04.1 K26.0 K56.69 R18.0 K91 K51.811 K28.1 787.6 542 568 777.50 K85.11 K12.0 K52.839 531.71 530.4 K90.2 K26.3 P76.0 K31.6 787.24 R19.32 K28.7 K05.329 532.3 533.90 534.10 579 K04.7 K22.11 K06.022 K08.199 K33 K59.00 787.22 K86.3 R19.33 R19.12 K14.2 K03.3 K81.9 K50.914 K44.9 K85.00 787.02 568.8 P78.84 560.30 K13.5 530.89 K85.12 K51.213 562.1 K13.0 K56.690 K01.1 K74 R19.34 R13.0 K83.2 K06.021 K28 K02.51 K75.4 K64.8 531.51 K95 K11 K11.6 P78.89 540.1 K38.2 K55.062 K51.018 569.49 533.1 K90.89 K51.814 K50.113 569.3 K91.61 K63.5 531.40 533.0 777.9 564.02 530.85 K05.20 K94.30 K43 P78.0 K51.50 K56.5 K05.5 534.31 K00.1 K56.2 R11.14 537.81 533.9 K04.5 K29.80 K31.4 K96 R19.03 532.10 K60.5 R19.15 530.84 K55.031 K31.3 K08.401 562.10 560.39 533.20 K12.33 R10.829 537.6 K08 534.21 K51.214 K82.3 532.21 K12.1 K66.8 K51.90 577.0 K29.31 K41.30 K76.1 R11.13 787.3 K56 530.7 R13.13 K57.52 K91.81 K02.53 K63.2 R10.10 K51.418 251.4 K02.62 K55.011 K62.82 K11.0 K91.871 K45.8 K70.41 K29.91 R10.825 R10.822 K44 K57.00 787.29 K20.0 K00.8 K50.813 K14.5 K51.00 K56.49 456.21 K04.6 K77 K80.45 558.4 K52.0 569.42 R19.5 R17 K57.91 537.1 K51.813 K08.26 K41.40 K05.01 K85 K76.6 K05.313 533.10 787.20 R10.819 K08.3 K25.1 K59.1 K08.89 K50.819 K68.19 R14.1 R11 K26.6 K51.019 K90.41 K94.33 K59.4 787.01 K20 R18 577 531.3 K65.9 K05.222 540.9 K03.5 K51.413 533.70 K59 K40.41 787.04 K91.858 K02.52 533.7 K50.012 K38.3 K41.90 K51.819 K56.51 K57.32 K50.019 777.1 K80.43 K22.8 K71.51 K08.113 K08.194 K64.2 R19.01 K92.1 I85.10 K76.0 K57.50 K80.60 K51.319 527.9 K57.81 K94 R13.10 K08.112 K03.6 K94.02 P76.8 K35.80 R10.83 530.81 532.01 K31 K71.7 K29.61 527.1 569.1 K74.0 569.82 K85.21 531.5 536.9 K52 532.20 577.9 K09.9 K57.10 533.00 532.60 I85.00 564.00 K50.80 531.6 K74.4 K31.7 K73.8 534.91 K28.3 K51 777.4 R10.12 K57.41 530.83 540 K05 K80.61 K41.91 534.1 K59.03 K81.0 533.40 K02.3 R19.2 K90 K58.2 K94.20 K08.132 K27.0 K80.33 562.0 K26.1 K95.81 K08.129 P76.2 K51.818 K11.7 K40.01 K82.4 K35 K43.5 K52.1 532.2 K71.6 K21.0 562 K71.3 K92.89 777.51 K63.0 560.81 K08.493 R19.35 K51.013 R15 787.1 777 533 537.8 533.41 K05.10 K13.3 534.6 K26.7 K25.2 534.50 531.50 K08.419 K22.710 K62.0 K65.4 R11.2 531.1 533.2 K15 K61.2 564.5 K13.24 534.11 K01.0 K55.052 K50.018 K94.22 K11.20 K13.21 K40.91 P77.1 K27.7 K80.47 K06 530.20 569.81 K66.1 K91.32 K50.811 K62 K12.31 K20.8 K08.439 K94.09 K24 558.9 K94.00 K56.609 533.71 564.0 534.61 K51.414 K51.012 K80.46 K23 Q43.1 K00.5 K55.059 562.02 K55 K05.229 K00.6 K06.1 K45.1 K55.1 K78 251 K31.9 K40.21 K08.491 K51.014 530.3 K31.819 K08.434 K80 K56.3 K51.519 532.7 K82.9 K70.40 532.70 532.5 560.32 527.5 K17 R16.2 560.8 527.0 K27.4 K72.11 K18 K80.62 K91.870 K08.402 540.1 530.0 K03.2 K91.850 K11.8 569.43 560.9 K13.22 K31.82 K08.131 531 537.9 K80.41 534.4 K50.10 K08.139 K43.9 K91.89 K58.0 K51.212 K85.01 531.21 K70.11 564.6 K40.90 K85.9 K03.7 K97 456.20 K56.60 534.41 K51.419 K06.020 K91.873 K05.311 543.0 532.1 527.7 K90.49 K74.5 K74.69 K02.9 K80.50 K59.3 K08.121 K46 K85.02 K92.2 787.03 K81.2 P77.9 K09 K05.322 K35.3 K61.4 K09.8 K31.2 K43.2 K26.4 560.89 K80.12 K42.0 K45.0 K50.011 569.84 K41.11 K91.71 K51.511 531.20 534.20 579.4 533.6 K08.191 536.3 K50.114 K11.21 K94.32 560.31 K75.1 R13.12 K56.7 K94.12 R16.0 K72.01 K28.6 777.52 K83.8 251.8 K90.81 K06.3 K14 K00.7 532.00 R10.814 K83.1 K08.56 787.99 K57 534.0 K80.18 K70.10 R19.8 K76.89 K08.24 R10.827 534.90 K06.2 K03.4 534.7 |
| Disrythmia | I49.8 427.0 427.89 427.9 I49.40 427.41 427.32 I24.8 I25.10 427.42 I49.49 I25.89 I20.8 I49.5 I49.01 I25.9 427.69 I25.811 427.5 I24.1 I25.2 427.2 I25.41 I25.810 I25.3 I21.11 427.81 427.1 427.31 I47.1 I21.29 I21.3 I47.9 I47.2 I24.0 I49.1 I49.02 I25.82 I21.4 I49.9 I25.42 427.60 I21.19 I20.0 427.61 I20.1 I21.09 I25.812 I25.83 |
| Dorsopathy | 721.6 M51.04 724.5 M54.6 M43.8X9 722.0 722.83 722.80 M41.00 I20.8 721.91 M50.00 M53.9 722.70 M46.80 M51.26 M50.90 I25.3 I21.11 721.0 M46.00 M46.30 M47.12 M48.00 724.1 721.41 M47.812 724.79 722.30 724.01 I25.42 M42.00 722.73 723.2 M46.47 M46.90 M51.06 M40.299 724.6 M48.9 722.92 M47.817 722.4 M54.89 M46.45 M43.6 724.71 I25.10 M48.20 M53.0 724.09 722.91 M40.00 724.00 M47.10 I25.811 M41.30 I25.2 M46.40 721.5 I21.29 724.8 M54.08 M48.02 M47.14 720.2 I21.4 M54.02 723.6 723.8 M41.40 721.1 720.81 M54.14 722.10 M40.40 724.2 M40.10 I25.83 720.9 M43.5X9 722.81 M40.50 723.1 724.70 M48.30 722.2 M45.9 722.32 722.39 M43.00 M48.04 M46.20 I25.9 M41.80 723.9 I24.1 720.1 I25.41 M53.1 M43.27 722.72 721.7 M51.9 722.90 M53.82 722.6 I24.0 M51.36 724.4 M54.14 M48.08 M51.44 720.0 M53.2X8 M54.30 I21.19 I20.0 M43.4 M47.819 722.51 I20.1 721.3 722.31 721.2 721.8 722.11 721.42 I24.8 M47.16 724.02 M51.46 M49.80 M53.3 I25.89 M54.5 M41.20 722.93 723.3 I25.810 723.0 722.71 723.5 724.3 M50.30 722.52 720.89 722.82 M50.20 I21.3 M54.12 M51.24 M48.06 723.7 M46.1 I25.82 721.90 M51.34 723.4 724.4 M54.2 M48.10 I21.09 I25.812 |
| Endocrine | 242.11 E23.6 253 255.4 259.8 253.2 E27.0 242.90 253.0 250.62 250.72 250.02 P72 250.80 252.02 242.1 258.9 E05.40 E20.9 P72.8 250.3 242.2 P74.3 P72.2 259.2 775.6 P70.2 362.03 244.2 H05.239 250.42 250.10 E26.02 259.50 E23.0 240.9 250.30 E24.9 E11.9 255.2 E11.359 250.82 250.50 250.70 255.9 259.51 255.13 242.41 250.0 246.3 E11.329 E34.52 362.05 P84 250.9 259.52 250.2 L68.0 E11.00 E26.01 E05.91 357.2 775.3 775.89 775.4 P74.5 241.0 253.4 250.52 E05.20 E11.21 P70.1 P73 E31.9 376.32 242.30 243 252.9 250.00 242 259.3 P70.4 E11.29 E34.9 362.02 250.22 775.5 E04.9 P74 242.80 E34.0 E05.30 250.7 E11.349 255.11 E00.9 E31.8 250.90 E26.9 250 E11.641 242.91 E11.65 E05.11 E22.9 E27.9 P74.1 250.60 P72.0 250.12 222.1 259.4 E21.1 250.6 242.31 240 E01.8 252.08 E04.1 241.1 259 E34.50 250.20 775.1 242.21 250.32 250.8 E07.1 255.41 255.10 362.04 255.6 E27.5 240.0 P74.6 E21.3 775.8 E23.2 E21.2 E11.40 244.9 242.20 E27.2 E23.7 E27.8 775.9 E05.41 252.8 242.4 E04.0 252.01 704.1 246.0 E07.9 775.0 E26.81 E11.51 E21.0 252 P71.8 E01.2 244 362.0 255.12 E13.42 P72.1 242.40 253.3 242.3 250.40 E03.2 E21.4 E25.8 250.5 E34.51 E11.620 775.81 255.42 E05.10 250.20 E04.2 255.5 E05.31 250.92 362.04 250.6 246.9 253.5 242.10 E21.5 E27.49 E05.90 E03.9 244.8 250.4 P74.4 255.14 258.8 246.8 E05.21 242.81 241.9 259.9 P74.8 E11.8 362.01 P74.2 E26.1 253.6 E11.01 E11.311 E34.3 241 P74.9 E11.339 255.8 362.07 E07.0 250.12 E34.8 E11.319 E22.0 E07.89 246.1 P72.9 253.1 E22.2 246 255 242.9 E11.39 253.9 E11.69 255.3 E34.2 253.8 |
| Esophagus | 530.5 K22.6 530.85 530.87 K22.2 K22.10 I24.8 I25.10 530.82 530.84 K22.4 K22.9 I25.89 I20.8 K22.70 K22.11 530.83 530.89 530.19 I25.9 530.12 530.13 K21.9 I25.811 K20.9 I24.1 I25.2 K20.8 I25.41 I25.810 I25.3 I21.11 I21.29 I21.3 I24.0 530.86 530.3 530.21 I25.82 I21.4 I25.42 530.4 K22.8 530.0 530.10 K22.5 530.11 530.7 I21.19 I20.0 530.20 530.81 530.9 K22.3 K21.0 I20.1 I21.09 I25.812 K20.0 530.6 K22.0 I25.83 |
| Health-Services | V58.49 V12.50 Z79.2 Z28.3 V04.5 V77.2 Z01.89 V06.3 V10.52 Z00.8 V12.49 V02.61 V70.8 Z93.52 V15.05 Z76.89 V03.89 V71.5 Z41.3 V03.2 V06.0 V02.60 V58.0 V15.1 Z22.31 Z11.59 V18.3 V15.02 V58.30 V72.9 Z82.49 V13.69 V53.7 V57.3 V03.6 V17.4 Z87.898 V12.01 V58.78 V77.0 Z02.1 Z71.3 Z48.01 V71.1 V06.1 Z22.51 V01.84 Z41.8 V54.19 V31.01 V64.2 Z38.1 S92.302 V64.00 Z51.11 V07.0 R13.10 V15.01 S59.101 V04.81 V59.9 Z09 V05.9 V01.79 Z02.81 V68.9 Z91.048 Z00.00 V80.2 V12.61 Z13.0 Z13.5 V67.59 V05.2 V09.0 V54.12 S92.302D V12.59 Z82.5 Z46.82 V54.89 Z91.012 Z11.8 V68.89 V17.49 V15.83 Z13.9 Z00.111 Z01.20 V02.2 Z72.820 Z22.50 V06.9 V76.3 V48.9 V72.7 V65.49 V67.4 V07.9 V74.8 V44.1 Z23 V54.09 V58.9 V55.4 Z91.011 V81.6 V77.99 V41.1 Z91.010 V61.29 Z91.038 V55.1 Z28.82 S72.471 S72.471D V29.3 Z87.19 Z48.812 V58.43 Z48.810 V70.2 V55.8 V81.5 Z76.2 Z48.813 Z47.89 V45.51 V41.0 V06.5 Z91.81 Z46.6 V13.02 Z97.3 V07.8 V71.2 V43.1 V75.7 Z13.220 Z01.110 V31.00 Z68.51 Z11.1 Z13.21 V58.74 V72.11 Z01.10 Z38.30 V82.5 V03.3 V41.6 Z45.2 Z11.6 V12.60 Z51.0 Z38.69 V13.7 V78.9 V71.7 V54.10 V54.9 Z11.2 Z55.9 Z39.1 V61.20 V44.52 V70.9 V33.01 V74.9 V05.8 V40.9 Z01.12 Z22.1 V75.9 Z53.29 V29.8 V49.75 Z43.1 S72.011 Z20.9 V54.01 S02.8xxD V20.31 Z13.4 V72.82 Z20.6 V20.32 F69 V12.09 V21.9 V05.4 Z96.1 V18.0 V30.00 Z86.79 V66.5 V04.82 V58.11 V67.9 V21.8 Z71.9 Z43.0 V62.3 V81.4 Z63.9 Z48.03 V67.2 V58.89 V15.09 V04.4 V78.3 V14.3 Z43.8 Z04.41 S42.101 V68.09 V53.6 V54.16 V20.1 Z98.89 Z16.11 V77.7 Z72.4 V27.4 Z83.2 V03.81 V53.2 V15.86 V81.2 Z46.9 Z13.228 V85.52 V12.51 V77.6 V54.13 V03.5 V55.0 V07.1 V72.85 V54.15 V67.00 V72.1 Z00.3 Z86.718 V15.5 V61.8 Z85.830 Z86.59 V73.89 P00.2 V72.19 Z46.2 V09.1 Z20.811 Z79.01 Z04.9 Z51.89 V49.2 S52.90x V79.9 V71.81 Z85.528 V82.6 V01.7 V44.0 V19.6 V04.6 V06.2 Z38.31 V45.89 V72.84 V85.51 V77.3 V34.01 V65.43 V65.40 V69.1 Z22.330 H57.9 V04.0 Z77.9 V72.6 V07.39 V01.5 V71.9 V39.01 V12.6 V72.83 V58.82 V78.0 V78.8 V58.31 Z68.53 V44.2 V03.82 V85.54 Z76.1 V20.0 V58.69 Z13.89 Z01.810 V72.2 V70.0 Z48.00 V70.5 Z13.83 Z01.812 V45.2 Z88.3 V70.4 V06.4 V58.73 V05.1 V58.61 V72.31 Z52.9 Z00.110 Z87.01 Z03.89 Z83.3 V19.8 V15.89 Z87.440 V02.4 V49.5 Z86.69 V72.0 S59.101D V01.9 Z20.3 V10.81 Z13.29 Z01.00 V25.01 V13.9 S52.90xD V50.2 Z38.00 Z93.0 V78.1 V02.59 Z03.6 V72.69 V57.89 V02.9 V58.62 V67.51 V74.1 Z68.54 V79.3 Z79.891 V68.1 V07.2 V30.2 Z93.2 V58.32 Z08 V69.4 Z43.4 Z87.798 Z87.09 Z91.018 Z71.1 V04.8 V04.2 Z86.11 V77.1 V76.12 Z22.8 V15.88 V29.0 V61.9 Z12.31 V82.9 Z38.01 V03.1 V02.0 V54.11 Z62.1 V20.1.411 Z63.8 Z48.89 S72.011D Z60.3 V16.3 V71.02 V71.89 V01.89 Z98.2 Z71.89 V14.1 V20.2 Z02.9 V58.81 V70.1 V18.19 V24.1 V65.5 V15.06 R68.89 V80.3 V79.8 Z51.81 V06.8 V72.5 Z89.519 V58.71 V30.01 V40.3 V71.4 V58.83 Z13.828 V19.1 Z12.6 V04.3 V19.2 V53.90 V29.9 V07.31 V12.00 V57.1 Z78.9 Z13.1 V77.91 Z68.52 V65.9 V72.62 V70.3 Z77.011 V14.0 V01.1 Z02.89 Z82.2 Z01.818 Z80.3 Z46.89 V05.3 V02.51 V72.63 Z04.71 V12.79 V17.5 V72.12 Z13.88 V05.0 V30.1 Z86.19 Z00.2 V15.03 V53.09 V72.60 Z76.0 Z30.011 V04.89 V03.9 Z04.3 V85.53 P00.9 V21.0 Z20.1 V58.3 Z84.89 Z93.1 V49.89 Z88.1 V40.1 V62.4 V21.2 V82.3 Z02.79 V72.81 S42.101D S02.8xx V41.2 Z48.02 V39.00 V67.09 V65.8 V65.3 Z20.89 Z22.0 Z11.9 Z04.6 Z46.1 Z00.129 Z16.10 Z83.49 V50.3 Z28.9 Z88.0 V71.09 Z13.6 P00.89 Z97.5 V02.1 V82.89 Z01.811 V57.21 V06.6 Z41.2 V64.05 |

| Category | Codes |
|---|---|
| Hematologic | D68.59 P59.29 D72.822 P58.42 P52.5 D50.0 P61.2 280.8 P57.0 P50.8 P52.0 282.69 286.0 774.5 D69.2 P50.0 P10.0 281.8 D57.411 D51.0 287.8 286.7 282.5 287.32 774.1 289.51 D55.9 P59 772.14 282.60 280.9 D53.8 D57.211 D59.9 284.8 D61.811 D59.2 D59.6 282.1 D73.81 P52 D61.09 287.3 281.0 282.42 773.4 P50.1 283.11 287.5 286.1 D68.318 D64 773.5 D68.8 285.21 776.4 772.3 D72.820 284.1 P51.8 D63.0 283 P61.5 D78.02 D59.1 772.9 P59.9 283.10 D50.9 D59 D68.311 D53 P10.8 D78.34 D75.81 D72.9 D57 P10 287.39 776.3 D58.9 P10.1 286.4 D74.0 287.31 D55.8 D75.0 285.29 D68.61 D78.81 D72.810 D69.9 D55.2 D60.0 D52 P55.0 D73.0 282.61 D59.5 P54.6 D53.0 285.1 285.22 289 D64.1 289.82 D59.0 D78.33 P61.0 D63 284.0 D61.1 D74.8 P54.2 289.81 281.9 D64.2 D76.3 776.7 P58.0 P58 D57.412 P58.41 773.0 772.13 D64.3 D61.9 D53.1 D68.52 D72.828 D75 P60 P52.8 D52.1 D60 P58.8 283.1 P57 D70.0 D78 D58.1 D69.49 289.59 P61.1 281.3 282.49 D68.32 283.9 D58.0 286.5 P52.4 D69.3 D73.9 D68.51 D57.80 774.31 D57.819 286.3 D78.21 284.81 D73.1 D77 282.68 P10.3 D56.4 D68.0 D73.89 D78.12 D58 284.2 D55 D76.2 P61.8 776.6 282.7 776.1 D53.9 289.52 D57.811 D56.2 D75.9 284.9 P52.6 287.30 D72.829 D68 772.6 287.0 D52.8 282.62 D68.69 776.0 774.4 D72.825 286.2 D57.812 P54.4 284.09 P61 D72.819 284.89 D73.2 282.40 D57.20 284.01 D68.2 D72.0 P58.9 D75.89 D55.3 D61.82 D56 P54.1 773 D53.2 282.9 D56.8 776.9 D69.41 772.4 776.6 D60.9 D78.11 P59.0 286.9 D78.32 281.1 285.8 D52.9 P50.3 287.41 D57.212 289.8 287.33 D72.823 D72.824 282.64 P54.3 P55.9 D59.3 D58.2 D60.8 D70.4 D51.9 D69.42 P59.8 P61.4 D64.4 D69.1 P57.8 D75.82 D69.6 287.1 P58.2 280.0 P58.5 P56.99 P61.6 772 P50.2 D56.1 P50.4 D70.3 P52.21 284.19 D78.01 282.2 D70.8 D51.1 D50.8 D61.818 D51.8 D57.02 D57.1 D75.1 287.2 774.0 P51.9 774 D69.51 P54.8 P58.1 282.41 285.2 D56.3 P54 D72.818 287.9 P55 P61.9 P55.8 772.8 D54 D52.0 773.3 282.6 D73 D61.2 772.0 774.2 284.12 D78.22 282.8 287.4 D64.9 D55.1 P59.3 P52.22 D74 P54.0 D70 280.1 287.49 287 D76.1 D78.89 289.7 P52.1 D72 D50.1 289.89 D72.1 P56.90 281.2 772.2 286.6 772.11 289.4 D68.1 D59.4 D69.0 D73.4 D58.8 P52.9 D67 P59.1 776.8 D59.8 289.9 283.9 D69.1 289.3 D61.3 D51.3 D61.89 D66.0 D57.40 283.0 D71 D56.0 281.4 D61.01 772.10 D69 284 P56 772.12 D51 774.6 D56.5 285 P10.4 D78.31 D63.1 D70.9 P54.5 D72.821 D57.219 P59.20 D57.3 D51.2 D64.89 P61.3 P50.9 772.5 D69.59 285.3 P54.9 D68.4 P55.1 285.9 776.2 289.50 D73.5 282.3 D62 282 776.5 D72.89 282.63 P57.9 D61 D57.00 D63.8 P51.0 D70.2 D74.9 P58.3 P52.3 D56.9 D55.0 D57.419 776 D68.62 D69.8 D68.312 D70.1 289.83 774.7 D64.81 773.2 774.39 D64.0 P50.5 P50 D53 D61.810 774.30 D76 P51 289.5 D73.3 773.1 D60.1 D68.9 P10.9 D65 P10.2 D50 D57.01 |
| Hyperplasia-Prostate | I24.8 I25.10 N40.0 N40.3 I25.89 I20.8 600.90 600.01 600.20 I25.9 I25.811 I24.1 I25.2 I25.41 I25.810 I25.3 I21.11 I21.29 600.00 I21.3 I24.0 600.3 I25.82 I21.4 I25.42 600.21 N40.2 600.11 I21.19 I20.0 600.10 600.91 I20.1 I21.09 I25.812 N40.1 I25.83 |
| Hypertension | 404.00 403.00 I15.9 403.11 404.0 402.00 404.10 I11.9 401. 403.0 I13.10 403.90 405.9 I16.1 402.10 402.1 404. 403.1 I15.2 404.11 404.03 403.9 404.1 405.09 404.91 I15.0 401.1 405.0 405.91 403.10 404.12 40311 404.9 404.13 I13.11 I12 I16.0 403 401 402.91 404.02 405.1 I10 I15.1 402.0 I15.8 405.99 405 I13.0 402. I12.9 402.01 I16 403. 404.93 401.9 I11 404 405.19 I12.0 404.01 403.01 403.91 405.01 405. 402 I15 I13 404.90 I14 402.90 404.92 402.9 402.11 I16.9 I13.2 404.04 405.11 401.0 I11.0 |
| Immune | H10.439 716.08 M86.642 273 003.23 M46.80 288.65 711.03 360.04 716.86 I89.9 I80.3 279.00 477.9 M02.369 K52.81 N34.0 711.11 582.81 728.82 K29.81 I89.0 715.15 597 710.0 A39.83 730.25 716.46 E10.11 M02.169 D73.81 M15.9 M81.0 J34.81 716.0 708.8 J32.9 M16.10 451.8 511.0 G37.9 M12.129 711.64 716.15 E06.0 H01.019 478.22 254.9 716.8 535.2 595.2 730.13 J45.991 692.84 J32.0 711.05 446.0 721.42 D89.0 M86.179 373.32 G35 D84.1 I05.8 555.9 571.49 716.43 250.23 D80.3 250.83 249.61 N03.5 250.21 535.10 459.10 H57.13 719.33 H16.309 733.01 476 451.84 249.20 711.90 464.2 D86.1 M12.829 373.90 D89.49 288.0 J35.03 493.81 J35.01 250.81 715.10 713.2 475.0 711.83 721.4 K51.40 711.86 M94.0 716.35 716.02 716.88 254.0 L41.8 289.3 364.0 D72.828 M12.00 535.71 711.58 M13.80 716.01 M05.60 692.3 H57.8 M02.339 H15.129 716.20 711.84 288.6 N30.30 M02.319 363.05 711.65 376.12 M86.149 373.31 716.6 396.1 M08.3 E32.1 H30.009 696.8 D86.82 D83.2 535.21 713.6 370.60 478.24 M12.80 695.15 363.12 582.9 M02.149 373.2 373.1 245.0 474.8 715.93 H30.029 I08.0 M13.849 E10.21 530.19 D81.0 464.10 L50.8 L24.0 M02.10 376.1 716.41 249.00 478.20 E32.8 I02.0 476.0 M12.529 711.95 464.20 N30.20 279.8 729.4 L40.8 716.38 D81.819 D84.9 E31.0 M12.9 715.95 M86.9 370.61 D72.819 692.6 363.22 M02.359 M13.829 372.10 686.0 360.1 396.0 249.5 535.4 711 D83.0 M01.X49 716.29 714.81 I77.6 H01.129 716.85 D86 E08.21 535.51 693.0 H30.129 397 711.21 D80.9 595.1 D72.824 446.3 364.11 242.01 M46.90 M13.149 730.01 M12.379 474.11 373 D81.818 I01.9 716.42 478.21 279.11 J34.1 L28.2 370.55 730.28 556.5 E08.40 711.52 373.3 595.3 249.90 582.1 730.07 I80.219 711.63 289.1 D80.5 M12.519 D81.2 288.02 K51.80 556.8 394.9 446 J37.1 288.66 364.00 715.30 D80.7 J44.0 711.76 730.06 M72.2 D89.811 D89.9 696.4 I87.009 M36.2 711.09 398.99 E10.65 595.81 446.29 716.91 711.91 711.18 E10.9 M12.579 715.04 H15.019 288.01 555.1 I87.099 694.6 M34.1 711.44 519.2 K58.9 477.8 364.1 D86.2 250.71 715.24 E08.65 711.24 M05.10 715.16 714.2 714.0 715.17 394 711.50 L20.89 716.84 464.21 M31.30 M86.18 464.01 716.49 G37.0 446.7 711.53 711.47 L27.0 M06.4 461.0 M12.88 J04.10 695.4 M12.19 H01.029 J04.11 372.31 H16.249 363.04 459.19 457.8 478.11 714.31 711.9 E09.65 M02.159 D86.89 493.11 M60.20 711.73 H00.039 J35.1 J45.990 716.60 730.20 711.60 D85 H16.429 457.9 L73.2 459.13 716.12 696.3 K29.01 M02.179 474 370 555.0 M31.0 711.43 363.14 711.08 H20.9 K29.41 245.4 H05.10 M33.20 363.03 379.01 716.54 250.41 373.34 L92.3 363.20 716.50 464.11 716.63 446.4 I73.1 694.60 511.89 715.98 249.31 461.9 L51.1 715.34 I06.0 M02.349 711.35 395 711.88 474.9 716.81 I06.1 711.15 D86.0 370.9 711.82 M13.139 D89.89 716.34 372.12 H40.40x0 M13.0 I07.2 I80.8 711.12 715.9 L04.9 716.44 730.00 730.21 J01.20 711.33 H30.139 692 D86.84 M02.18 535.50 M00.09 H20.819 708.1 M19.079 M12.319 716.51 556.2 M13.159 L10.9 288.04 289.6 M35.00 279.09 730.0 E08.10 695.12 363.07 459.1 370.63 716.87 692.9 L30.4 364.04 M13.869 451.89 715.21 392.0 M12.119 720.89 249.4 686.09 711.20 H30.93 711.40 715.26 K52.2 L21.9 716.95 D89.2 D80.2 J91.8 379.93 M12.339 J30.5 693.8 715.80 M17.5 D83.1 720.0 249.71 461.2 571.40 711.5 L92.9 D81.7 358.0 H20.13 730.24 249.60 396.8 711.81 M01.X8 580.4 695.11 363.1 716.96 715.97 250.53 E08.620 475 716.53 249.70 M08.40 J35.9 D80.4 D70.0 478.29 715.27 M02.129 G70.00 H16.269 451 M00.08 373.9 L08.89 693.9 710.1 686.01 K29.90 H10.409 M30.0 451.19 396.2 363.00 530.10 M12.169 M47.10 M14.80 714.3 493.1 E10.39 716.99 M13.88 M02.00 E05.01 477 K29.00 E09.311 N03.8 E04.1 698.2 H16.449 728.71 J30.0 D80.6 464.0 711.29 461 288.3 571.41 L44.8 721.41 719.36 597.89 474.00 N00.9 711.98 I08.8 711.89 M12.50 D72.825 713.0 564.1 721.91 711.87 711.80 H20.829 N08 H44.019 716.22 341.1 I06.2 715.36 443.1 E10.29 D72.0 279.05 715.13 J01.40 457.0 373.00 711.45 579 M47.16 I01.8 H20.23 277.39 M01.X39 716.65 L93.0 289.2 461.1 535.41 279.2 M19.229 M17.10 I88.0 E10.51 H40.40x 730.26 J01.90 M00.019 M19.219 715.18 D70.4 716.11 708.0 695.89 711.28 250.43 M86.669 716.16 M12.18 K75.4 M02.379 L53.0 711.31 249.7 H30.039 H10.429 288.61 394.1 D86.9 H20.049 556.9 694.8 D70.3 D70.8 L50.3 M02.39 730.04 451.0 D87 325 288.60 473.9 M02.19 I09.89 694 E10.40 D84.0 698.4 E06.9 H05.119 250.31 L25.9 M00.049 711.69 692.8 L52 245.3 706.3 711.32 K51.50 E06.4 720 M08.00 711.94 396 K29.80 M86.659 716.06 D76.1 D82.9 711.49 279.03 580.81 715.06 397.9 535.61 535.00 715.08 L25.1 M31.4 715.90 288.51 358.00 K51.90 716.93 J45.909 D82.1 716.55 370.6 H01.9 J39.2 713.1 I80.00 M12.369 360.00 J04.0 D84 696.0 730.23 250.03 H15.059 373.12 K29.91 288.2 M86.68 D70.9 J32.1 715.35 I88.9 250.61 473.3 379.06 711.54 715 M19.249 M12.39 I80.10 692.1 730.08 394.2 716.48 M47.14 364.23 K51.00 711.48 716.17 373.11 719.30 370.62 M19.279 446.20 535.5 099.3 716.9 711.68 M30.3 372.1 711.72 716.07 379.0 J35.2 249.6 715.23 716 364.22 H01.009 719.3 716.89 716.81 711.70 730.14 720.2 708.9 582.0 N30.10 715.89 556.4 580.89 398.90 N03.3 716.21 711.14 398.9 L51.3 D89.82 H44.029 398 E05.00 716.67 719.37 451.83 711.42 710.3 D86.83 363.2 396.3 H16.339 360.01 711.7 249.9 370.54 399 493.8 711.6 582.4 H15.009 715.31 245 714.8 249.30 E32.0 451.2 D83.8 535.70 D81.1 I87.029 E08.9 245.8 511.8 714.4 719.39 696.1 279.12 711.37 711.16 715.2 M19.039 M16.7 597.0 M46.1 L23.0 372.13 457.2 M15.0 711.06 692.4 H16.409 370.50 715.25 J45.21 I80.209 I88.1 713.7 713 708.2 555.2 M31.1 395.2 462 M12.869 242.00 J01.10 H01.8 288.64 K29.61 373.01 379.00 730.27 E08.51 M00.069 L51.9 711.2 288.5 473 379.04 493.0 715.09 714.32 E32.9 457 J94.1 721.1 249 535 493.11 M30.5 D86.85 E08.8 373.33 J05.0 373.0 716.40 D81.5 493.12 715.96 692.81 K50.80 493.00 D81.4 J30.2 H01.119 K73.8 M19.029 279.3 J02.9 595 N03.9 I08.9 L50.1 H00.019 J30.1 695.13 L50.5 716.27 M16.9 711.17 730.19 H20.019 364.10 L98.1 595.0 L44.0 L30.1 530.12 716.0 288.69 730.22 E06.1 556.0 711.22 469.0 A02.23 683 364.24 730.18 279.02 451.81 H44.009 716.09 J39.8 N30.81 D89.3 254.1 373.13 730.09 K52.1 720.81 364.02 M14.60 288.62 464 M12.149 716.30 M19.049 I80.9 391.8 715.33 L00 711.36 571.42 358.01 511.9 364.21 711.56 279.13 711.71 446.5 D83.9 711.66 M00.029 711.02 D72.829 370.52 711.78 L88 363.11 719.31 M00.079 360.02 695.1 249.21 D81.89 711.13 M12.859 715.05 D80.0 571.4 M12.179 D88 M15.1 M13.10 N00.3 M19.90 713.3 H30.119 372.14 K20.8 735.2 705.81 D81.6 H16.439 464.00 M35.3 711.3 464 719.32 478 363.21 250.73 J01.30 250.01 M31.6 245.1 478.8 692.0 694.9 G70.01 693.1 493.82 D81.810 L27.9 715.91 714.9 M13.819 M12.329 J32.4 477.0 597.8 I05.1 478.19 I97.2 714.89 E08.01 363.13 711.0 493 720.8 364.05 474.01 M01.X69 691.8 288.50 493.01 289.53 459.12 H16.329 249.91 H30.23 D82.2 716.90 535.40 I89.1 |
| Infections-Bacterial | 041.4 T79.A21 T79.A29 041.86 A43.0 034.1 020.1 030.0 041.5 A27.89 038.2 023.3 T50.A25D A25 021 A24 040.89 A48.0 036.1 041.7 040.82 041.10 039.3 T50.A24S A39.83 T79.A22A 036.3 A23.9 B96.81 A44 041.04 A48.51 026.1 A37.80 032.84 041.00 A41.02 H90.A21 A22.0 A40.B B95.0 A27 A21.2 A49.3 T50.A22 T79.A29S A48 020.3 040.3 041.84 038 A30.3 020.5 T50.A22D I82.A23 H90.A32 A39.50 T50.A21A A02.2 A36.2 T50.A23A A41.89 038.19 041.81 A20 A30.9 020 A22.2 A41.2 A32.89 041.03 030.8 A23.0 036.81 A28.2 A36.85 A39.9 032.2 041.6 A49.8 I82.A29 036.42 B47.9 M79.A21 T50.A25A 020.8 M79.A22 038.40 041.83 A31 038.44 T79.A22S 040.81 027.2 027 032.83 A46 038.49 039 041.85 036.40 A44.0 T50.A23A 038.42 A42.89 B95.7 A21.8 A41.01 A36.9 A39.3 A20.9 A40 H90.A22 T50.A26A A49.01 A37 023.1 041.09 A49.9 A49.1 A22 A23.1 T79.A21A A48.52 A49.2 A22.7 A42.7 A28.1 M79.A29 034 A32 031.0 038.10 A38.9 A39.2 037 A21.3 A48.3 038.0 A48.8 A39.1 022.0 A49 020.9 A28.9 A39.89 026.0 036.0 036.82 A30 T50.A26 A44.1 041 T79.A21S 036 A26.7 A28.8 039.1 A22.1 022 B95.1 A41.1 A36.0 A44.9 038.3 020.0 A27.0 038.12 A20.0 A43.8 023.9 A23.8 A39.53 027.9 B96.4 T50.A23 A44.0 040.2 A25.0 A41.51 A32.12 A36.86 A23.2 A39 022.8 029 B96.89 A22.9 030.1 034.0 036.43 B95.8 A38.0 T50.A24 B95.4 B96.5 T50.A26D 027.8 A40.8 T79.A29D A36.1 A39.82 A43 A44.8 040 039.2 T79.A22 021.9 A43.9 L08.1 A26.0 020.2 A37.10 036.41 022.1 041.9 A42.82 038.41 A41.50 A35 B96.7 021.1 T50.A21S A49.02 T50.A21D A32.7 T50.A23S A20.1 T79.A21D A37.11 A48.4 T50.A25 A36.81 A41.53 033.8 A30.8 032.82 A21.1 A36.83 032.81 A21.0 A39.2 A32.11 023.0 A40.9 041.19 M60.009 021.2 032.1 035 A40.3 B96.1 028 A36.82 022.3 A34 A37.01 036.2 A47 023 B95.2 A38 A30.1 033.9 A42.2 041.81 032.89 A30.4 A39.4 A39.0 I82.A22 020.4 A41.4 027.0 038.8 A39.52 021.3 036.89 T79.A22D A21.7 A32.9 031.2 A37.81 041.89 A37.91 T50.A23D A30.2 A23.3 T79.A29A A42.9 030.3 A32.82 A41.1 A36.8 T50.A26S 038.43 B96.3 H90.A31 A31.8 039.4 041.01 032.9 039.9 024 040.42 040.41 030 A39.51 041.82 041.2 A36 030.2 A41.59 A48.1 A26 T50.A25S A38.1 025 021.8 A31.1 A37.00 040.0 026.9 031.9 027.1 036.9 A42.0 032 031.1 023.8 A24.0 A28.0 A43.1 B95.61 B95.62 A20.3 041.11 A31.0 I82.A21 A25.9 022.9 A42 A41.3 B96.6 B95.3 A30.0 041.12 038.9 A45 A28 A30.5 A48.2 031.8 A20.8 A31.9 A39.81 T50.A24A 040.1 A31.2 A39.84 A25.1 A42.81 033 T50.A24D A24.1 A27.9 A26.9 A27.81 B95.5 A23 K90.81 021.0 041.05 A37.90 032.3 039.0 A36.84 A42.1 033.0 A33 032.0 026 A32.0 A29 T50.A22S 023.2 032.85 A24.9 A21 A41 T50.A21 039.8 A24.2 038.11 A20.2 A24.3 A20.7 041.02 A41.52 022.2 A36.9 A38.8 031 |

| | |
|---|---|
| **Infections-Fungal-and-Other** | B46.4 122 B90.2 M86.642 B76.8 B58.81 M90.829 123.0 B87.2 B77 T50.B96D B69 136.4 B63 B37.1 B54 117.4 122.3 M62.3 B40.81 B67 B47.1 126.2 B60.11 B94.2 M62.89 122.4 B95 114.4 B35.0 133.8 138 123.1 115.00 B73.1 B74.8 B81.2 128.8 728.82 B57.1 728.3 B57.49 B48.8 B60.12 B35.2 139.0 116.2 B65.8 B37.81 730.91 T50.B93S 117.7 B65.3 B39.5 B94.0 B88.2 M90.869 B67.39 115.94 730.00 127.3 730.21 B75 728.85 B51.0 125.2 730.25 B37.3 B96.81 B40.89 134.9 133.9 M90.839 131.01 131 B43.1 B83.9 B57.32 B65.0 136.29 B36 B37 122.2 B87.89 B95.0 B50.0 B66.8 114.9 B66.2 730.81 B43 B44.89 M90.849 115.10 B38.89 B40.9 M62.838 B71.1 B83.2 130.1 730.94 B42.89 112 B60.0 B35.5 B74.4 B38.7 B52.9 M86.639 B98 117.8 B97.32 137.2 B42 B97.11 127.9 B50 730.27 B37.6 122.5 B58.9 B81.3 115.03 I32 B39.1 B37.41 122.6 B58.89 730.98 B56.1 115.99 B46.9 730.13 112.9 B55.9 T50.B96 730.95 B59 A07.8 T50.B91 T50.B92D B46.0 B67.4 M46.30 123.6 B76 B77.0 112.89 B37.84 M86.179 136.5 115.14 B73.00 B47.9 B36.0 B60.10 B45.2 730.77 730.36 111.0 B38.2 B80 120 730.38 B40.3 B48.3 130.9 B40.2 121.5 B37.0 730.82 730.32 B83.1 112.1 122.8 M89.659 M62.81 110.4 T50.B94 B53 B43.0 730.03 117.5 111.3 B81.1 T50.B91D B46 M86.619 728.10 127.8 128.9 B88.8 130.2 M86.19 730.29 137.3 136.3 730.24 B97.0 113 B57 136.9 B87 B95.7 728.81 B97.89 B45.8 B45.9 110.6 121.2 B76.1 M46.20 B37.9 B97.10 B74.3 B57.0 A59.00 B66.1 B77.89 112.5 730.19 115.13 B37.83 111.1 728.0 B55.1 B35.1 B78.9 117.0 730.22 B85.2 H32 M86.60 B62.00 M86.20 730.17 112.84 730.18 B38.0 B71.9 728.12 B53.8 B79 B60.2 M89.619 T50.B93 B97.35 M86.149 B45.1 B37.5 A59.9 730.09 B38.1 B67.7 126.3 M90.879 728.6 B51.9 111.8 B35.6 123.4 M61.10 B78.0 121.1 B58 B44.0 B40.0 B72 B97.4 B35.4 730.76 730.15 730.39 T50.B95A B37.7 112.83 730.75 B96.22 117.2 M89.669 B36.8 128.71 B41.7 B71.0 B71.8 A59.8 B48.0 110.9 B56.9 B65.9 M35.7 B67.99 127.7 730.96 B44.9 T50.B95D M62.82 B65.1 B88.3 L94.6 730.88 123 B68.1 B46.2 M60.10 B68.9 B95.1 B86 B43.2 T50.B96S B87.81 112.81 114.0 728.2 B60.19 131.9 B57.42 120.9 131.8 B97.39 B83 M89.679 137.4 125 130.4 730.33 B96.4 117.1 B78.1 B99 B97.5 133.0 133 130.0 121.4 B88 T50.B92S 116 B57.2 B37.2 A59.03 123.5 B85.1 B78 B41.8 M86.9 B38 112.2 139.1 B97.29 B96.89 B51 123.9 B38.4 129 B40.1 B97.12 123.2 B40.7 728.5 B99.9 B88.9 M61.40 B71 B81 114 730.31 J17 121 B36.2 B50.8 B95.8 B81.0 132.9 B57.39 M72.4 T50.B91A B55.2 B69.81 134.1 132.0 B39.2 T50.B93D 110 B42.1 134.8 B77.81 126.0 B67.5 B95.4 T50.B91S B41.9 B96.5 B48.4 B69.89 B67.61 132.2 B46.5 M24.20 B85 B35 115.04 T50.B92A 123.3 B58.00 B37.42 B47 115.05 B68.0 110.0 B36.3 B85.0 B87.9 115.01 730.26 126.8 110.5 T50.B95S 121.6 730.30 730.01 B56.0 130.8 B60 B87.82 B83.4 728.19 730 B51.8 728.11 B67.69 111 B90.9 B94.1 B96.7 M86.669 B52.8 B96.21 728.89 730.12 B38.81 B45 B73.02 B67.8 B58.2 124 730.28 B67.32 121.0 110.8 B87.3 134 D86.9 132.1 A59.09 M90.819 111.2 730.07 B66.4 115 125.7 130.3 B39.9 B85.3 B67.2 728.84 730.04 B69.1 A59.01 134.0 B38.9 B40 B94.9 115.12 B55.0 B42.7 127 B36.1 B97.33 B93 136.8 M60.009 128.1 B74.2 B96.23 121.8 B85.4 M90.80 A59.02 730.10 112.4 B39.3 117.9 B96.1 730.06 M72.2 B73 B82.9 B66.5 B69.0 139.8 B95.2 137.1 B42.0 B46.3 B48.1 115.09 B67.90 B94.8 137.0 M86.169 130.5 M86.139 131.09 115.02 T50.B93A M61.00 125.1 120.1 128.0 139 M89.69 728.87 M86.659 730.16 132 B78.7 B35.3 M86.129 116.1 120.0 B44.1 122.0 B39.0 B45.3 730.37 730.72 B57.40 730.85 B97.21 B82 B74.0 M86.119 128 M86.69 B39 B94 136.21 B58.83 728.88 B57.5 B39.4 T50.B94A T50.B92 B38.3 120.8 B99.8 M86.159 B96.82 730.78 B74 728.13 B48 730.80 B82.0 B36.9 B46.1 T50.B94D B68 B37.82 B96.3 B43.9 728 B58.3 122.9 728.79 T50.B96A 730.35 730.11 B62 B90 B97.81 115.91 M89.629 B81.8 115.19 B96.20 B42.81 B90.8 B97.6 B84 M72.0 132.3 126 B97.34 B41 B65.2 127.4 B45.7 730.02 M90.859 128.9 B82 146.9 128.9 B89.619 T50.B93 M89.619 B50.9 B52 M86.68 B65 114.3 T50.B94S B52.0 B96.6 B97.30 B66 B89 B87.4 B95.3 120.2 M62.10 B97.19 728.86 B37.89 B57.41 B58.09 136.0 122.7 130 136 127.1 B90.1 730.71 G02 B92 114.2 B87.0 M60.20 730.84 730.86 730.99 112.85 B64 114.1 730.92 B61 730.20 M89.649 125.9 115.93 B42.82 B73.01 127.0 136.1 117.6 730.05 730.73 115.11 121.3 B57.31 M86.679 B97.7 B35.8 M86.629 730.14 B76.0 B87.1 B83.0 B73.09 B96.29 M90.89 127.5 B83.8 B60.13 121.9 123.8 B95.5 B44 B53.0 126.1 B58.1 112.3 B41.0 B97 730.74 728.83 B46.8 B47.0 B37.49 125.3 B91 B96.0 130.7 131.00 B58.01 B74.1 B56 B53.1 730.70 127.6 B83.3 B66.3 B67.0 B66.9 139 112.0 120.3 B44.2 B55 B58.82 730.79 B42.9 730.93 B48.2 728.4 115.15 131.03 122.1 730.97 112.82 137 118 115.92 125.4 B67.31 B88.0 M89.68 B35.9 B96 M90.88 115.95 |
| **Infections-General** | 041.4 017.10 B02.21 055 098.14 042.9 B09 011.80 569.5 B69 136.4 A19 078.0 057.8 A50.09 017.81 590.3 083.8 016.4 013.83 003.8 B67 341.21 372.05 016.32 380.12 015.62 017.20 466.1 380.2 001.1 001 A52.10 B94.2 H70.009 B35.0 A92.8 017.44 017.82 123.1 115.00 016.71 B74.8 T79.A12S 081 016.60 011.8 B02.24 T50.A15A A05 B60.12 T50.A25D A04.9 B33.23 T50.A11A 682.6 098.41 B35.3 116.2 N37 018.42 040.89 A77.0 117.7 095.8 045.1 A07.4 A63 B08.1 018.01 070.3 B88.2 M90.869 482.42 A85.1 005.9 041.10 015.5 375.30 A39.83 016.36 013.53 B96.81 011.93 488.11 041.04 A48.51 B26.1 A73 A08.32 B43.1 A18.4 B65.0 A41.02 T80.A19D 018.85 B16.9 066 B87.89 060.1 I82.B29 324.9 P36.2 012.85 685 B50.0 A27 A66.5 070.70 A66.1 B02.9 T50.A93A P35.1 A48 100 B17.9 590.1 058.1 079.3 130.1 B06.00 B42.89 112 017.9 B60.0 T50.A22D B74.4 B52.9 046.2 077.9 381.1 086.0 483 056.01 A54.41 016.24 062.4 A83.9 K67 B97.11 449 085.0 007.8 013.62 A92.9 P36.39 T50.A94S 016.1 011.5 J05.10 466.0 A50.07 B00.82 015.61 B81.3 017.76 012.81 B08.010 015.04 771.3 A30.9 063.8 016.93 A52.05 B02.0 A22.2 059.22 041.03 008.69 011.23 381.0 A06.9 380.1 083.1 A28.2 A36.85 B46.0 013.14 070.52 B77.0 082.40 054.1 482.89 112.89 A02.0 T50.B12S 010.02 032.2 A18.54 B27.82 B73.00 008.5 047 036.42 T80.A10D M79.A21 B45.2 A08.39 681.0 111.0 T50.A25A B40.3 B48.3 375.31 094.1 017.33 014 J31.0 038.4 G00.3 B83.1 045.23 011.43 045.01 T50.B94 041.83 071 015.21 B43.0 016.53 048 094.8 472.1 016.04 010.94 073.7 A96.9 017.4 079.52 015.12 127.8 016.14 T80.A11S A00.1 097.9 A07 011.64 A66.3 038.49 039 100.81 137.3 136.3 382.2 B57 006.0 038.42 A51.45 B87 B95.7 771.1 H70.209 483.8 012.1 A36.9 010.01 063.9 121.20 B97.89 420.91 B26.85 B45.8 081.2 074.22 A88.0 014.06 042.0 567.3 115.13 B55.1 016.54 016.35 B35.1 573.3 053.13 A37 015.75 079.2 J18.9 N11.8 008.45 056 A52.06 A54.83 A23.1 093.81 A53.0 424.90 A54.84 I40.0 H10.239 A48.52 A49.2 112.84 013.23 016.41 099.55 051.0 004.9 093.20 B14 B24 373.4 011.72 B79 098.39 M89.619 T50.B93 A80.9 091.1 A32 A17.82 031.0 M90.879 043.1 488.1 A51.9 J09.X3 A87.2 123.4 G00.9 A21.3 B40.0 B97.4 066.49 A02.24 038.0 B06.01 A48.8 017.06 730.76 421.1 098.5 094.83 020.9 001.9 115.9 015.66 A02.22 A50.03 B37.7 026.0 077.98 482.8 730.75 117.2 055.0 A93 A77.8 A04.5 A01.05 B02.32 003.21 A53.9 322.1 B48.0 B56.9 A44.1 098.50 015.70 B17.2 B67.99 H60.399 098.59 T50.A94 T50.B95D G03.9 A01.2 123 039.1 077.4 682.2 P35.2 T50.B15D B31 017.66 A80.39 A81.09 381.3 B87.81 B60.19 B32 038.12 074 321.2 131.9 012.8 A43.8 B57.42 131.8 015.92 482.2 A07.0 T50.A13D 125 011.90 072.2 060.9 014.05 027.9 101 B30.9 A06.1 320.2 B97.5 484.8 A41.9 130.0 013.90 078.19 T50.B12 084.8 040.2 102.2 A25.0 B57.2 123.5 I31.2 J15.5 B27.81 099.59 B08.72 M86.9 005.2 094.82 H59.42 112.2 T79.A11 139.1 J85.1 003.9 008.43 003.0 A82.9 480.3 A08.2 083.2 123.9 070.43 078.4 078.1 B38.4 B40.1 B00.89 034.0 A51.5 012.01 B71 006.4 513 J17 074.3 A56.4 A18.12 A77.49 B95.8 132.9 B55.2 011.15 017 008.42 T50.B93D L03.039 B77.81 126.0 A56.09 077.3 B95.4 484.5 323.71 682 074.0 008.49 B17.11 010.05 A66.8 A85.2 B46.9 B96.5 A51.31 054.42 011.41 A52.8 132.2 A50.45 A04.7 A77.9 B46.5 006.5 115.04 009 A39.82 B37.42 011.0 B08.09 054.73 A44.8 012.84 P36 B85.0 383.21 059.0 040 126.8 T50.B95S 771.0 074.23 323.51 685.0 A26.3 J28.0 B01.2 B83.4 A42.82 B00.7 B67.69 015.7 038.41 B94.1 421.0 021.1 098.86 A68 T50.B15A A49.02 T50.A21D 013.22 A32.7 B58.2 124 099.49 B67.32 079.83 121.0 A54.5 A52.3 J12.0 B87.3 070.1 B08.03 132.1 016.55 111.2 T50.A25 090.7 A36.81 A41.53 018.80 A81.01 B69.1 A21.1 098.82 I30.9 002 102.9 016.94 J04.30 A19.8 050.0 017.65 A18.18 A50.43 B93 023.0 A92.5 103 B74.2 032.1 121.8 B85.4 A18.59 A59.02 015.22 A40.3 392 T50.B50 T80.A19S B82.9 A52.15 T50.A93 A30.1 B42.0 053.11 376.00 B48.1 018.04 137.0 033.9 130.5 045.9 131.09 A05.1 A19.9 A77.1 A02.9 T50.A11 125.1 590.01 011.46 A82 128.0 A39.4 139 M99 016.0 041.0 070.0 008.04 B18.1 015.73 B78.7 B34.0 P36.10 B44.1 B39.0 077 095 383.22 A93.0 381.4 062.3 036.89 326 B82 046.72 422.91 099.0 070.30 A54.85 B08.79 016.31 372.03 057.0 013.51 A37.81 T50.B94A 017.73 464.30 T79.A29A A50.9 A54.03 A52.71 730.78 A30.2 J31.2 A02.20 010.8 372.30 A26.8 012.33 008.67 B96.3 B43.9 018.84 A18.52 B58.3 I31.4 A79.89 046.8 T50.B96A B10.89 039.4 B97.81 B05.3 091.8 686.1 115.19 B81.8 B42.81 T50.A95S I38 A69.20 013.96 A06.81 015.2 B05.81 126 070.20 010.84 A93.1 041.82 A56.3 A65 085.9 054.40 093.89 015.54 016.42 084.7 464.31 B05.0 111.9 030.2 A48.1 008.64 072.1 H62.40 B10.82 A06.89 055.1 052 J12.81 466.11 T50.B95 007.2 A38.1 060 A31.1 B34.3 372.04 026.9 M72.6 016.12 027.1 590.2 P39.2 126.9 110.2 011.35 095.1 480.9 013.13 B44.81 004.8 A20.3 013.65 013.8 052.8 I82.A21 382.00 B76.9 A18.10 013.54 A68.0 M89.60 A56.19 116.0 041.43 A52.09 320.82 B25.1 J12.1 A17 036.8 091.3 M89.639 001.0 074.1 A41.3 A03.8 B02.7 B43.8 B97.30 018.96 682.5 682.0 372.20 B87.4 072.9 018.81 J16.0 G03.0 A83.0 041.12 J15.6 321.3 B37.89 012.16 019.2 B57.41 B58.09 070.5 K65.2 045.03 A98.8 070 095.4 G02 A66.9 T50.B14D 053.29 B33.0 A48.2 031.8 P39.0 099.41 A06.6 A31.9 114.2 J03.90 065.1 730.99 H05.029 A55 A50.7 H67.9 047.8 017.3 054.41 050.2 127.0 016 B34 018.94 011.34 033 056.0 115.11 121.3 A42.81 A24.1 079.50 B76.0 T50.A93D 072.71 099.51 015.95 A15.9 012.10 059.12 B07.0 T50.A13S B96.29 017.34 B00.2 B60.13 P35.9 078.2 049.8 121.9 B44 598 103.0 A05.0 H04.309 021.0 041.8 B22 B02.31 730.74 014.81 012.83 372.21 011 B91 A77.2 B58.01 J15.9 T50.B13D 015.51 054.49 B33.3 026 A54.33 A29 053.21 A51.39 682.4 488.09 321 T50.A12S T50.A22S 088.0 032.85 A00.9 008.8 073.8 730.79 A80.4 B42.9 A41 A71.9 A80.1 093.23 567 039.8 A06.0 H04.429 P37.0 015.23 078 112.82 041.02 115.92 A41.52 B88.0 B96 022.2 A36.89 A52.9 B35.9 I31.8 A38.8 013.92 I82.B23 H66.90 015.94 M35.8 G06.1 391 B46.4 017.86 A02.25 A17.83 122 B90.2 B19.11 T50.A96A T80.A11D A43.0 A74.81 016.51 056.00 A18.39 472.0 102 099.5 A67.0 H70.229 323.4 J09.X9 117.4 015.71 A80.0 B47.1 041.5 126.2 J21.8 484.7 I33.0 016.13 B95 016.95 054.71 043.2 422.93 045.10 128.8 B57.1 059.10 B37.49 091.4 038.2 016.61 021 567.31 B65.8 079.81 B37.81 730.91 018.86 B65.3 372.0 B02.23 B39.5 A48.0 B18.8 482.0 381.10 005.3 017.6 464.50 A50.6 P37.8 T50.A96 A52.11 B67.39 H04.339 T50.A14A 039.3 421 T50.A24S B16.1 B51.0 I09.2 053.19 015.56 T79.A22A T80.A19 A52 G05.4 051.2 A44 B33 079 A67.1 380.10 T79.A12D B27 053.22 103.2 771.5 A37.80 B27.19 I82.B11 018.8 A22.0 A18.83 A40.0 |
| **Infections-Respiratory** | 017.10 013.95 017.86 A17.83 J12.9 013.42 016.51 011.80 J15.1 017.83 017.42 A18.39 A19 010.91 J09.X9 017.81 015.71 017.30 013.83 J03.00 T50.A16D A15.8 016.32 017.20 015.62 017.31 B25.0 J20.6 482.32 J21.8 017.02 484.7 016.13 012.31 016.95 017.44 017.82 016.71 013.61 016.60 T50.A15A 011.33 011.84 018.93 016.61 011.10 485 015.50 018.92 018.86 J01.21 J12.3 017.70 482.0 018.01 464.50 011.31 J11.82 482.42 010.06 014.01 J01.20 016.61 015.55 T80.A19 013.53 010.14 A19.0 011.93 013.44 A18.4 A17.9 017.22 T50.A15D J20.9 J03 017.80 T80.A19D 011.91 018.85 A18.83 013.06 016.03 012.85 A17.89 J18.1 012.21 015.01 482.82 011.42 462 017.32 J20.0 482.31 017.74 015.14 012.00 J01.10 487.1 012.06 483 012.86 017.93 J09.X1 011.30 016.24 J02.0 J15.212 016.73 010.11 016.66 J10.81 013.62 016.30 480.0 A18.6 013.10 016.34 017.25 017.13 016.23 012.04 J01.00 016.45 015.11 013.14 482.89 010.02 A18.54 013.93 J10.08 J01.41 011.96 J07 011.26 015.10 J11.08 018.05 466 J05.0 461.8 017.33 014 015.55 011.43 011.01 011.05 016.53 015.20 012.03 011.92 010.00 016.93 A19.2 013.80 012 A15 016.05 488 A18.85 015.53 011.23 012.04 J01.00 016.45 015.11 013.14 482.89 010.02 A18.54 013.93 J10.08 J01.41 011.96 J07 011.26 015.10 J11.08 018.05 466 J05.0 461.8 017.33 014 015.55 011.43 011.01 011.05 016.53 015.20 012.03 011.92 011.00 016.96 482.1 J11.1 482.84 011.24 011.04 J15.3 483.8 010.01 010.83 J15.20 J20.1 A15.0 017.63 J22 016.90 013.26 A19.1 010.16 014.06 016.92 016.35 016.54 010.96 J00 015.75 016.91 J18.9 T50.A16 012.35 016.01 015.76 017.92 A18.31 010.93 012.13 018.95 016.21 T79.A19 013.23 016.41 012.14 482.49 014.85 011.72 017.61 J47.0 J09.X3 T50.A16A 013.66 011.53 013.55 015.83 A17.81 J01.80 015.74 J12.89 B44.0 J18.8 464 482.83 017.06 A18.89 016.65 011.14 J10.01 015.66 J11.2 017.90 013.31 012.34 018.91 J08 010.80 J20.5 A16 I82.A19 011.76 017.24 011.50 013.20 016.16 016.50 015.03 015.65 012.32 015.70 016.52 A15.5 461 016.44 A18.14 015.52 A22.1 011.82 A18.53 460 J01.11 486 017.66 480.8 464.10 J02.8 012.30 J10.89 J21.9 J85.3 012.23 016.76 017.84 016.02 465.09 015.92 482.2 014.84 A17.0 017.40 484.8 013.45 464.20 013.90 014.00 J11.83 013.11 011.54 482.41 017.56 J15.5 012.22 482.30 J18.0 015.90 013.12 011.20 T79.A19S 013.81 013.33 J85.1 480.3 J01.40 015.85 012.02 013.43 016.74 010.95 017.64 461.3 011.95 012.01 016.46 A15.6 513 J18 016.85 J17 013.21 A18.12 011.17 011.15 015.84 017.13 011.90 016.63 014.04 013.40 011.94 A18 J01.91 484.5 J06.9 J20.4 464.00 461.1 010.05 012.36 010.86 A18.50 J03.80 011.41 015.15 J01.30 016.33 483.1 T79.A19D 013.32 012.84 012.05 017.52 011.06 J01.90 010.12 015 013.15 J14 011.45 J15.211 J09.X2 016.20 015.02 013.04 J01.81 014.86 011.05 013.82 487.8 015.16 011.66 A18.15 J10.82 017.53 010.90 M79.A19 017.43 013.22 012.11 J03.91 J12.0 J21.1 484.6 A18.2 017.15 J01 011.00 013.03 017.16 016.55 J15 016.70 018.80 488.0 013.41 013.25 017.04 016.94 011.32 013.64 J04.30 J09 465.9 A19.8 017.65 A18.18 J03.01 480.2 A18.51 016.72 J15.29 016.40 010.13 016.00 A18.59 015.22 016.75 018.00 017.51 513.1 T80.A19S 013.86 015.00 A18.82 017.50 018.04 018.06 013.85 J11.61 015.55 J19.2 J20.2 015.06 011.22 011.46 011.73 J06 016.22 010.03 013.30 017.60 015.73 017.05 012.25 017.46 016.10 T50.A15S 012.82 J04.2 016.31 013.51 A18.81 013.36 J12.17 017.73 017.91 464.30 018.03 010.10 A37.91 016.43 013.16 010 481 017.35 A18.17 015.86 011.56 J20.8 T50.A16S 011.52 017.72 010.00 011.36 015.82 012.33 011.03 483.0 A15.4 018.84 013.35 A18.52 011.51 012.24 018 J11.89 J03.81 J15.0 013.96 G40.A19 010.84 017.95 010.92 015.54 464.21 016.42 464.31 026.85 A48.1 011.63 464.01 464.4 012.12 T50.A15 J04.0 513.0 015.00 J20 011.11 484 466.11 011.25 013 J12.81 017.23 015.96 016.12 013.63 484.3 461.0 017.03 011.35 480.9 J04.10 013.13 017.12 013.65 015.63 015.80 015.93 A18.10 A18.7 T79.A19A 013.54 J04.11 482 016.62 011.75 A17 J12.1 016.56 J10.83 017.36 011.13 J10.1 011.85 016.64 018.96 011.61 015.91 015.64 482.81 018.81 J16.0 013.05 J20.7 011.86 015.81 J15.6 017.71 012.16 A15.7 011.43 016.37 015.13 J02 018.90 013.81 J15.4 017.14 482.39 A18.32 482.9 A18.13 015.24 465.8 011.60 010.82 487 017.96 013.94 012.26 013.34 016 017.55 018.94 011.34 011.21 014.82 018.02 015.95 J20.3 017.26 017.00 A15.9 012.10 017.54 J11.81 013.56 015.72 013.02 017.34 016.15 J10.2 012.15 011.74 480.1 J05.11 011.44 017.45 014.83 015.13 017.11 014.84 A18.01 010.81 466.19 J11 015.25 A18.84 016.06 016.23 017.21 014.81 A18.16 012.83 J85.2 014.80 017.75 011.16 484.1 011 017.61 J15.9 017.61 012.12 463 J05 J04.31 A18.03 J85 016.96 482.40 012.80 013.24 013.84 J04 464.11 011.62 J19 017.85 014.02 J21.0 015.23 011.81 T80.A19A J13 013.92 461.9 015.94 011.02 465.0 |

| | |
|---|---|
| **Injuries** | S06.9X4S T84.490D T63 S63.621D S82.232C S65.391S T85.614 S25.809D T85.398 S72.346 S83.011A S60.212A S66.526D T80.59X S52.232D S53.026A S82.115P S56.822 T86.391D T63.614D T50.992S S61.221S S19.9XXA 943.23 T23.179 S71.131 T80.410S S93.129 T40.0X2 T81.83XS S72.436R S63.433D T24.591D 832.01 S72.91XK T24.031D S36.533A S56.196S T22.232A T24.332D S92.355K S59.129 S31.112 T78.3XX S66.120A T81.516 S21.442D S82.043P 803.50 S63.232S T53.0X4A S92.123A S36.020D S72.019A T63.511S S82.864R T65.831S S63.657D S14.157 T04 T84.620 S22.040K S56.892D S72.441A T56.1X1D T83.85X T53.7X3 S68.123A T47.1X6S T26.71XA S14.4xx S89.002D T53.1X1 S62.307S T65.213S S82.242M T22.741A S62.023D T46.6X1 S32.131K S52.515C T65.0X3S S46.012A S52.265C S92.221B S65.111A S06.1X1A S60.461 S61.549D S36 814.03 S24.151A S50.811A T62.1X1S S42.446P S52.264G T24.491A S52.042A S92.599K S06.4X9A S92.416D T20.411A T81.507D S12.600B S82.492P S21.351S T45.1X4S S32.000 T75.4XX T82.391 T43.623 S36.290 T54.3X1S S42.209R T56.3X3D S59.222 S99.829S S09.19XD S34.22XD S63.652A 942.45 S32.471A T40.3X2 S60.444D S82.146S S82.822K S50.10X T83.022S S56.128S S05.91XD S82.874Q S42.111D S25.402A T82.322 S62.614S S86.929S S90.464D S74.00xA T65.4X2 T61.784D S83.30X S32.456B T53.93X T24.501A S37.69XA S37.091D S63.022D S80.849S T46.991S S61.021 T49.0X5 S90.411D S30.826 T63.834 S64.494A S42.033D S32.058D S42.191D S52.332C S31.110S S82.871M T43.4X5S S95.911S 861.20 S20.309A S01.01X S13.171D S72.366M S82.154B T79.5XXA T83.69XS S09.90x S00.252S S83.096S S01.90XA S72.031K S63.439A S06.365 T50.Z92 T25.139 S09.312A T50.Z11D S37.031 S82.426A S82.022K S95.809 S12.300K S35.338S S92.244D T23.799 T28.412 S82.254 S62.312 T22.022S S02.118B S06.6X2D S60.940S S31.20xA T17.408D S09.312D S52.399F S62.251K T63.071D S82.464Q S52.243J S82.465Q S82.401N T49.8X6A 914.3 S55.212 S59.102K T86.23 S96.899 S22.050S S41.112D S50.11XA S52.272M S72.416C S82.876G S22.039D S42.144G S63.227A S66.518A S72.416S T85.511A T48.4X4 S82.034R T20.69XD S82.62XP S42.112B S52.244 S92.031K S12.151K S42.216 S41.031A S15.202D T23.079D S39.91XS T19.4xxA S82.291A S92.491K S01.449S S52.361J T63.813S T63.713A S72.434J S82.016D S99.029B S22.002A S80.919D S83.231D S82.209H S90.933 T78.07XA T53.2X1S S20.461 T36.7X2 927.8 T82.897A S72.365F T20.72x S63.065 S72.111B 821.00 S63.432D S65.911D S96.102A 944.46 S30.857A S83.001 S22.032D S87 S93.324D S15.9XXS S41.031D T63.93XD S82.154 S60.031 T23.469A S21.429A S52.282 T71.223 S62.338G S82.443D S67.90XS S92.111P 887.7 S15.121D 917.7 S66.516S S70.12XS S23.143A 956.0 T39.8X3S S72.012D S90.01XD S61.205D S82.242S S82.843M T46.6X6D S39.002A S62.331B S09.312 T45.2X1S S82.843D S86.111S S20.02XA S42.335A S59.299A S20.112S S81.019 S32.120S S83.123A S92.331G S50.912 S50.02XA S67.194 S02.630K S82.033P S66.002A S61.234S S61.217D 901.1 810.12 S89.392P S60.342 S52.265N S60.471 S00.451D S53.091S S72.041D S59.101G 923.03 T24.031A S01.81X S89.101K S51.842A T52.2X4 S92.323D T50.4X5S S82.234J T22.351D S76.311A S52.044R S82.872M S83.262S T55.1X2 T71.234A S82.156K T38.4X2D S31.145D S62.399P S63.217A T41.0X2 S92.411D S63.121D S52.225S S56.329D T23.091A 866.10 S45.011A S36.99XA S72.361A S52.343F 949.4 S75.811S S81.822D S82.124P S72.442D S46.891 S92.326D T25.219A S72.425R S42.251B S66.120S S88.021S S27.69XD S52.513J S60.569D T63.832A S65.508S 935.1 S82.113N S62.329D T83.39XA T63.823A T61.04XA T81.532 S42.035 S02.602A S52.609A S56.418D 806.31 S82.253Q S98.912S T21.31XS T37.5X6S S09.0XX S72.364D S25.192S S25.511A S59.919S S62.637B 812.19 T33.519A S00.221 S01.401 883.1 T24.501 T18.120D T14 T48.0X5D S40.862D S30.870S S92.242G S48.021A S52.035S S52.321B S52.334R S62.620P S42.409A S43.402S S62.513D S82.126A S92.051D T43.226A S62.637A T21.33xA S62.644P T48.1X6D S45.392 S52.223A T45.1X4D S92.356K S72.052K S52.365C T65.6X1A S72.341H T82.827D 844.8 S89.312K T37.8X4 S34.101A T84.121A S10.92XD S61.344D T20.63X S92.513 T23.321A S62.226P S12.040S S93.409S S03.02XA S45.901D T39.1X1A T81.519D S62.102S S63.391S S37.022A S62.601 T43.601 S15.209S S46.109D S95.999S S52.381D S83.242A T42.6X1S T17.890D T43.605S S82.136S S72.102M T51.1X2D S04.041A S43.025A S82.102Q S65.409S S92.243G S22.49XA T22.362S T53.4X4 S61.111D S11.11X S32.402G S42.123S S62.300G S36.92XD S62.609A T53.92XD S20.169 S42.223 S82.391C S68.619 T71.21X T63.414S T52.8X3S S85.142D S81.841 T44.1X2 T82.838 S82.391M S61.253 S04.10X S93.302 S92.422 S52.202C S62.211K S98.922D S37.051A S72.302P S92.226P S92.101G S95.902S T45.2X6D S62.031D S26.92XS S72.144 945.16 S65.508A S75.911A S49.031S S09.399 S70.922 883.0 S15.199A S99.191K S82.436B S37.899 S72.302B S04.041D S21.149S S28.211D S52.379S S62.145K S32.615A S41.121A S80.822A S20.109A S92.061P T21.60XD T36.93XD T24.239A T40.1X2A T83.86XA S01.82X T65.823 T65.94X S33.101S T85.868A S92.251B S02.670S S96.029 S95.101S S29.099S S60.322D S77.10XD S92.136P 863.93 S92.334B S63.228D S82.142R S52.246K T20.619D T26.31XA T25.391D T24.512S S82.423 S90.934 S82.461K S92.238D T63.866J T53.3X1A S11.95X T51.2X4 S62.337P S70.241 T85.122S S62.92XG 915.4 S62.311A T63.634S S72.322B S60.559D S82.192P T41.205D S89.099P S07 T65.214A S52.511G S36.30XA S65.091D S66.313 S01.422D T36.6X2S T50.991 S75.999D S92.401P S26.91x S85.131 T87.0X9 S52.326N T17.910S T50.4X1 S65.999A T49.7X4S S55.811 S32.509G S62.164D S91.141A T23.479S S72.90XS S62.335S S93.333S S43.91XA S52.263F 913.1 S82.862S S56.122S S76.011 T83.25XA S31.122D S62.186B 943.31 S92.351B S10.14X S45.36 S70.369A S31.639 S89.199G T56.1X2D S75.199S T52.4X2A S31.551S S52.132B S62.308K S62.629P S61.541 S89.312P S72.451A S82.022F S93.602D S66.510 T36.4X3D T23.442A S32.442 S37.818S S66.911D S92.343S S30.842S S99.129G S63.093 S12.291K S32.601G S68.711D T45.695 S27.53XD S22.068K S92.233D T80.A9XA S90.02XD S00.04XD S62.224K S82.442K S04.032A S42.482P T49.0X1S S75.091S S56.019A S52.026N S82.846A S52.261K S68.114A S20.222A S92.513B T41.41XA S52.611H S14.107D S92.141S S52.336F S32.008B S12.112S S42.331G S63.251S S82.436Q S30.846 S52.512Q S82.453M S52.002P T23.442S S01.339 S06.9X4 T63.023S T25.412A T36.5X6 S42.265B S53.193S T63.114S S37.591D S92.919B S31.210S S92.522G S24.149 S42.125S S99.222G T48.204D S60.473S S09.21XD S62.638A S92.812K T40.904D S82.043M T20.32X S12.391D T21.17XD T22.791D T83.420S S61.243S S62.176S T23.449 S43.121S S62.025P S83.252D S00.01X T22.052D T79.A12D S30.93XD S25.892S S82.422A S20.462D S45.019A S42.463K T63.061A S27.812S T51.3X4 S69.90XS S93.135S S32.444S T23.179A S99.142S S25.502 S56.012D S81.839D S82.092J T50.2X1D S50.322A S00.471S S32.811K T85.860 S61.353A T82.221S T34.3XXD S84.20XD S72.133M T54.2X2 S09.399A S62.251S S76.822S S82.141B S78.911 T83.23XA S93.303 S34.103D S52.599P S21.232 896.3 S06.6X3S T45.7X6 T38.4X4 S82.235D T25.321A T36.8X2S S42.115G S14.135D S72.324R S82.846D S96.999D S63.635 T83.83XD T50.993S S06.4X1D 848.40 S65.811A S66.401S S82.312G S52.311D T16.2XXA S82.102H S52.246C S56.802S S06.384S T65.291A S36.269A S02.621A S83.211A S72.326M T17.598A S36.299 T23.672S S42.156 S52.132C S63.423A S82.431K S66.501A S62.252S S20.342A S42.492K T85.613 S32.409K S01.342A S60.542 S63.312A T67.5XX S51.822A S66.902S T33.90X 805.02 S92.035D T33.02X S83.122D S65.291S S30.0xxA |
| **Integumentary** | L10.81 L89.813 709.2 L02.639 R20.8 L02.619 L97.218 L92.8 L21 690.12 L56.2 Q84.0 778.1 521.15 L64.8 L03.123 L06 R23.9 782.8 L94.9 L03.213 L16 L97.503 L89.814 L98 L97.923 L66.3 L97.812 L02.422 L89.514 L02.33 680.9 690.18 L97.509 J70.1 680.0 L97.824 L89.149 757.5 L40.4 L97.413 L51 L63.0 L92.2 L97.504 L93 L24.89 704.00 P81.0 529.2 L02.432 709.01 K08.51 P83.8 L08 L89.124 525.0 L89.129 K05.00 L10.2 L72.2 K08.403 L89.614 K13.29 L03.212 L71.1 522.7 703.8 707.02 L49.7 L02.235 690.8 695.3 705.0 737.33 L89.209 L94.1 L04.2 729 K03.9 L02.239 L89.115 523 L29.1 707 L98.422 L89.620 778.3 782.1 L89.610 L41.8 778.0 L40.3 L35 521.09 710.9 L43.2 L03.114 L12.2 L64.0 L81.0 L49.5 L97.205 L97.306 525.64 L26 521.32 L89.103 L92 K08.530 L97.929 528.6 L89.629 L66.1 Q84.3 L89.603 L90.2 L12.0 K08.0 K05.30 695.57 L08.81 L70.3 L97.314 K08.9 L04.8 L50.8 L24.0 L89.503 L20 558.1 M35.9 692.71 L32 L89.009 L58.0 525.42 L97.118 L74.2 L89.40 L97.309 L30 L70 L40.8 L98.498 L69 525.52 L81.5 525.44 L89.501 P81.9 521.00 L63.6 K08.52 K02.7 L97.516 528.00 L95.9 L89.020 L23.1 L11 698.9 695.53 L40.53 L03.039 M79.0 L29.2 521.49 K04.3 697.9 L97.424 L71.8 L68.1 680.2 L40.1 L89.011 L24.9 L76.21 K05.21 522.9 L40 525 K05.32 K14.1 L97.321 L89.613 L97.404 K04.90 L82.0 700 K13.23 L03.891 L89.102 L28.2 L01.1 704.2 L89.601 L14 R23.3 521.11 L74.8 L73.9 L81.3 L97.904 L49.1 L85 L91.0 521.13 707.03 L97.209 L44.2 L55.9 L89.300 L81 P83.1 692.7 697.8 L97.121 523.30 L97.106 703 L74.510 521.9 R22.9 L02.214 757.3 L02.435 L94.0 L02.532 L03.316 L41.0 M96.2 695.51 L97.512 L20.89 L89.139 K14.3 L62 525.10 508.0 L02.519 521.02 L97.221 L89.304 L98.493 L20.9 L27.0 522.6 L56.3 L89.42 523.1 L56.9 L21.8 K08.104 690.11 778.2 L98.9 523.41 L03.011 L75.1 L89.110 523.3 L89.120 709.09 L03.323 L85.8 L49.2 522.0 L97.909 525.50 707.19 525.51 L97.304 L47 L89.524 L97.326 M79.9 L61 A69.0 L72.3 704.3 L89.143 701.2 L97.325 L03.019 K14.5 L02.838 L97.419 K04.6 L97.109 523.23 L98.0 729.7 525.3 L41.3 L44 K52.0 L03.324 705.2 L03.012 R17 528.1 L97.101 L10.1 729.72 L02.232 L57.5 L89.140 K05.01 L13 P83.6 L48 523.5 L97.305 L97.906 L01 L89.142 707.23 L97.928 L23 L97.522 523.11 L51.3 K08.3 L90.3 694.3 L75.8 L64.9 L44.1 L03.126 L97.325 L85.0 L76.11 L97.324 706 529.4 L02.811 L23.6 L03.327 L97.204 K03.5 523.42 L10 L97.813 L70.8 L03.041 L89.010 M79.89 L97.228 L94.3 695.5 655.63 L97.525 757.8 L97.403 L98.3 L89.214 L89.522 L89.201 782.5 L97.222 521.0 L98.423 L23.0 L93.1 L01.09 L97.904 529.1 705 L89.003 L02.431 L91.8 P83.5 L03.211 L84 L29.9 521.8 L29.0 L67.8 L02.225 L90.0 L89.621 525.61 782.9 K03.6 L74.511 525.7 L04.0 525.72 L02.411 782.4 R22.43 707.20 L97.114 Q82.1 L23.2 L89.311 R22.42 L98.414 702.11 707.8 528.72 L51.9 L68.2 L97.312 L89.302 L89.322 529.3 529.8 L53.9 695.58 L97.224 692.70 L40.2 525.73 L98.495 L97.913 L11.8 L89.222 704.0 L97.123 L31 L97.119 L20.0 L89.600 L53.9 L81.4 L50.1 P83.39 L37 L57.4 707.9 L50.5 R60.9 K02.3 L40.51 L89.150 L98.1 521.21 L94.4 L44.0 L57 L30.1 L02.531 L89.819 L98.419 L91 L98.425 L97.125 L54 L82 L89.45 M79.A29 L41.3 L03.115 522.4 L30.0 L89.811 L02.828 702.0 L02.32 709.1 L20.81 L95.8 698.8 L02.222 525.19 K05.31 L00 L89.314 680.5 523.6 508.1 R22 692.72 L03.12 L89.230 782 L94.6 L24.01 L10.00 692.82 L02.224 K08.419 L97.802 L98.428 L03.129 L02.226 729.99 L88 L03.029 L97.901 L57.8 L97.903 701.8 L02.433 L13.1 L92.0 L23.4 L03.122 L97.908 729.73 L97.915 K13.21 525.13 L89.210 L53 L89.113 L02.03 L97.316 707.12 L97.412 695.9 L97.215 L97.104 690.10 L25.8 L81.1 L03.119 L12 K08.439 L97.811 L03.91 L49.0 L03.111 L03.031 L03.898 L90.1 L02.216 L27.2 L77 L03.121 L02.611 L04.1 L02.213 L89.93 L97.103 705.21 L66.0 L97.414 L97.416 L44.3 L08.1 692.77 655.61 L63.1 L90 L41 L89.619 L27.9 L24.5 701.0 706.0 521.5 M79.A19 L02.12 L08.9 L97.924 L66 L02.425 L97.229 L89.023 522.5 521.01 729.79 L58 L97.808 L90.4 729.71 694.1 L24.81 L74.4 680.3 L97.809 L41.4 L89.104 L97.301 L97.208 L97.315 L89.202 L02.821 757.9 L05.02 L53.3 L97.816 O35.6xx0 L08.0 L04.3 704.02 L60 L89.013 K13.22 L89.504 L98.2 L05.92 L49.3 528.7 L02.412 L89.512 K08.402 L97.223 525.63 L74.513 K03.2 707.00 L53.8 L97.816 O35.6xx0 L08.0 L04.3 704.02 L60 L89.013 K13.22 L89.504 L98.2 L05.92 L49.3 528.7 L02.412 Q82.4 L80 L97.313 L97.329 521.4 L89.523 L02.211 L55.1 L49 L23.7 L94.2 K03.7 707.2 R20.1 529.6 L89.151 L89.131 L56.0 L02.414 L13.0 L55.0 L89.004 K02.9 |

| | |
|---|---|
| **Ischemic** | 410.72 410.51 410.41 I20.8 410.20 414.12 410.00 414.00 I25.3 I21.11 410.40 I25.42 414.10 413.9 410.12 414.2 410.32 410.61 414.01 I25.10 411.0 414.19 410.22 410.60 410.02 410.92 414.02 I25.811 I25.2 I21.29 410.30 I21.4 410.31 410.21 413.0 I25.83 414.8 410.81 410.82 411.89 414.07 410.90 410.11 I25.84 410.50 I25.9 414.9 I24.1 I25.41 410.80 414.06 410.70 414.04 410.01 410.62 I24.0 410.71 I21.19 I20.0 I20.1 414.03 I23.0 412 410.42 I24.8 I25.89 414.11 I25.810 414.3 411.1 I21.3 I25.82 410.10 410.52 414.05 413.1 410.91 I21.09 I25.812 411.81 |
| **Metabolic** | 269 E83.30 270.5 P71.4 276.7 264.1 712.38 M1A.9xx0 274.82 E71.42 278.02 274.19 278.01 276.69 263.8 712.86 E72.8 268.1 275.40 274.03 712.28 M10.30 E66.2 783.1 E74.9 E71.41 270.8 712.19 E83.59 E87.70 M11.80 E50.7 E88.1 275.42 E67.0 260 M10.9 E46 274 E72.20 M11.88 276.3 712.84 E83.50 P19.0 264.3 E50.8 278.2 275.1 E50.1 M11.20 E83.81 R63.4 775.9 277.81 268 712.36 277.88 E83.40 E53.9 271.9 P71.3 E87.1 M11.879 270.7 E52 D84.1 273.9 266.2 264.4 265.1 274.8 712.1 M11.249 712.8 M11.849 277.82 E72.9 M11.269 712.16 E50.4 712.88 M10.40 M11.28 269.9 E88.9 264.9 268.0 E55.9 278.0 712.85 266.9 277.9 E75.21 712.3 712.89 712.27 712.97 277.8 P72.8 M11.859 278.3 264.0 278.03 G93.9 E53.0 E16.1 271.3 251.0 775.89 E70.0 P71.1 E51.11 P70.1 275.5 P70.4 277.7 272.3 712.15 E50.6 277.87 712.39 274.81 712.96 712.10 274.02 E87.8 712.90 E15 273.8 R63.0 271.2 E66.9 277.86 266.0 E53.8 E64.3 263.1 278 E87.0 272.5 E66.3 712.32 E67.1 262 712.92 M11.9 M11.29 712.31 272.1 P19.1 P71.8 M11.279 712.95 277.85 E65 E88.40 775.81 783.22 E61.4 783.2 263.2 269.0 261 E83.89 P71 M11.829 P74.0 E74.12 E83.52 M83.9 269.2 269.1 783.21 712.33 M1A.00x E87.4 E87.5 P70.3 268.2 271.8 E87.3 M11.259 E76.01 E87.2 266.1 P71.2 E88.01 E55.0 251.2 266 273.4 274.00 272.2 275.01 275.8 712.98 E56.8 E56.1 E74.21 E44.1 272.7 712.91 263.9 E88.81 E78.89 E78.9 712.34 269.3 274.01 263.0 P70.9 712.22 277.2 E74.39 263 E67.3 276.6 E88.09 E50.3 E72.10 R63.5 E43 P19 265 275.9 264.5 E66.01 E83.00 E80.0 E88.89 712.18 E83.118 275.4 E63.8 264 274.10 277.5 275.3 712.94 712.12 E44.0 E40 M11.819 276.2 274.0 E51.8 270.6 E79.8 274.1 267 E83.9 712.23 276.9 272.4 P71.0 G93.89 712.2 712.87 E83.51 P70.2 271.1 276.4 E45 712.93 277.89 P84 P71.9 M11.89 E87.6 E41 712.17 265.0 E78.2 272.9 277.6 274.9 276.0 E56.9 E70.40 E83.110 M11.229 268.9 712.25 272 270 330.2 271 712.35 M11.839 264.2 275.09 712.30 276.1 R63.6 264.8 277.1 270.9 P70 712.21 278.8 272.6 712.9 275.03 M10.00 270.3 270.1 E88.3 E50.0 N20.0 775.8 P70.0 E63.9 E71.50 278.00 712.29 E74.4 274.11 E83.10 712.24 275.2 E78.1 P70.8 712.99 330.3 712.82 712.81 M11.219 M1A.9xx E50.9 251.1 E78.5 E78.6 272.8 P19.2 E71.318 264.7 E71.0 783.0 712.26 712.20 269.8 M1A.00x1 P72.9 265.2 M11.239 712.14 275.41 278.1 |
| **Musculoskeletal** | 553.9 M84.574D Q67.5 M84.442S M70.30 M25.761 M46.87 550.00 M60.231 M71.829 M85.622 M84.675A M84.753D D08.429 733.20 M84.472 M60.239 M84.346K 727.6 M12.332 735.1 M84.550K M62.241 M71.062 M66.219 M05.051 M24.275 552 M67.479 755.67 M89.322 M27.59 M25.842 M80.821S M84.753A M84.573P M54.02 718.80 M53.86 M85.58 Q65.00 733.43 717.41 M24.174 718.9 M85.012 M48.58XG M15.9 M10.452 M41.22 M20.21 M84.522K M16.10 M23.041 M66.211 M19 M13.822 M24.231 717.82 M90.871 M87.035 M48.57XD M96.679 M21.40 M32.12 M93.841 718.82 M46.96 M65.341 M11.10 M25.232 M67.279 M80.019S 719.99 M66.172 M84.462K M84.442K M96.639 M80.011S D08.231 K43.7 D08.849 M23.007 M12.212 M48.02 M84.619G M65.171 M67.229 M66.272 M60.819 719.90 M11.049 M84.553A M85.812 M89.312 719.70 M06.811 M12.449 M45.0 M79.4 M26.03 M24.075 M86.151 P13 M24.676 M99.77 M84.651 M10.329 M89.752 M23.91 M05.271 755.30 M80.052 M05.842 D08.829 M96.662 M00.232 M80.069A M85.311 M63 M94.279 D08.839 718.47 M77.30 M80.012G M05.372 M15.4 M62.511 719.04 M89.542 M84.561S M84.572 M94.9 736.05 551.01 D08.469 M94.0 M71.129 M80.011P M99.9 D07.649 M67.832 M05.519 M25.642 M84.533D M54.15 M05.069 M62.259 M84.662K M24.031 M20.001 M90.551 M60.031 M05.60 M61.146 M24.651 M41.112 M93.811 719.12 M89.619 M08.911 M14.879 M08.3 M05.69 M00.162 M05.722 756.55 M05.431 M18.32 M70.951 M80.859G 730.39 M92.12 M48.43X M65.852 M26.00 M23.052 M24.572 M40.14 M80.079D M86.542 M05.049 M65.261 M65.011 M84.663P M60.172 M05.811 M11.822 M71.562 M85.321 M99.21 M67.971 733.1 M84.364P M02.10 M60.074 M83.1 M35.9 M97.42XS M12.529 M05.879 M87.032 M46.38 754.5 M90.672 M50.23 719.50 M85.859 737.21 M12.9 M84.452S M84.443 M11.869 719.45 M93.879 M36.3 719.25 M89.49 M76.899 D08.929 M71.869 M89.042 M70.21 M87.334 M10.451 M26.23 M85.00 M06.831 727 M80.831 M90.852 M76.42 M89.28 M65.332 M86.661 M23.631 M50.823 M75.00 M05.479 M84.379 M10.321 M86.021 M61.129 M79.0 M84.569K M61.229 Q79.4 M08 M62.441 M66.811 M06.859 M54 M45.1 M84.571G M84.672D M23.005 M47.13 M84.612S M84.542 M67.272 K08.23 M12.379 M24.839 M10.019 M89.439 M48.51XS M89.30 M25.149 M22.3X2 M89.061 719.20 M84.756D 552.3 M60.831 717.0 M89.08 M32.11 M90.861 M42.09 M54.13 M12.519 M92.291 M93.851 719.83 M23.212 M71.80 M87.838 M06.059 M84.462A M99.43 M66.842 M23.322 722.30 M84.663K M01.X51 M07.60 M24.376 M80.069P M20.031 M72.2 767.5 M61.161 M02.361 M00.811 M26.213 M87.861 M84.673 M85.071 M84.475A M84.519P M99.18 M71.552 M24.575 M43.04 M12.579 717.8 M89.78 M12.219 M84.639 M84.563P M89.165 M60.261 M80.00X 730.37 M34.1 K40.10 P11.3 M26.82 M84.322G M84.564 M46.52 M51.05 718.89 M88.88 M05.121 M05.10 M84.674D M84.342D M89.571 M80.872A M84.841 M23.352 M71.051 M93.969 M99.0 M63.869 M48.061 M02.122 M21.379 Q65.30 719.92 M21.221 M25.361 M87.29 M84.359P M90.511 M72.0 M11.032 M85.639 M60.262 M80.00XP M84.350A 524.07 M94.261 718.41 M86.432 M54.81 M84.459 M10.422 M84.446K M84.572S M23.009 M06.4 M72.6 M87.263 729.39 M14.659 M87.274 M41.41 M84.750A M84.639A M84.672S 719.49 M91.10 M99.00 722.72 M14.621 M61.59 M67.942 M80.869G M89.032 M24.562 359.5 M60.019 M76.61 M24.371 M87.850 M92.219 M84.2 M26.52 M84.442 M84.334S M84.576P M13.112 M10.112 M86.252 M94.1 727.3 M85.342 M85.879 524.06 M85.821 M10.359 M84.475S M12.532 717.7 M27.62 M26.73 M43.01 M12.861 M60.869 M84.459A M17.5 M85.061 M48.55XD M88.839 M76.892 M19.022 M80.861A M84.321S M65.311 M80.852P M96.843 524.1 M80.019G M93.839 737.34 M84.632P M90.811 M18.11 M43.20 M24.673 717.81 M14.629 M61.032 M05.461 M85.052 M24.642 M67.89 M89.163 721 M35 M87.20 M99.41 M05.169 M84.575K M08.40 732 M05.349 M42.9 M08.80 M84.434A M86.559 M12.852 M70.11 M51.25 M94.221 M96.5 M25.28 M66.351 M84.345S M84.752S 526.69 D07.611 M87.261 M70.969 719.17 738.7 M84.312P M33.93 M66.371 M84.662A M48.54XS M89.124 Q65.1 M86.039 M65.821 M25.075 718.0 M84.532A M67.841 M75.52 M80.00XD M17.31 M11.142 M66.311 M01.X22 M84.652S M05.062 M84.652S M80.019A M89.162 M60.046 M32.9 M84.631D M87.132 M23.8X2 M84.552P 722.6 M43.4 M87.271 M84.542D Q67.7 M05.631 726.7 M62.58 M84.751 719.77 M84.671P M96.65 M99.31 M26.06 M62.529 M62.28 M12.121 M60.032 736.22 M71.862 M79.609 M87.00 M05.49 M48.34 M25.039 M24.232 M92.9 M20.42 M80.859P M90.579 M62.459 M20.039 M99.33 733.82 754.44 736.29 M06.321 M84.753P M60.279 M80.861 M84.433D M26.24 M33.02 M80.029G M11.221 M84.443A M65.121 M06.011 M25.531 M21.239 553.2 727.0 M66.852 M86.469 M70.869 M01.X39 M50.221 M84.662S M99.47 M80.052D M84.839 754.42 M10.072 M87.135 M31.31 M87.076 M19.229 M85.842 738.9 M12.361 M79.3 M48 M05.222 M48.56XD 719.56 M84.353A M77.00 M80.851S M80.051D M24.271 735.4 M84.362G M86.569 M80.862K M63.842 526 M80.069D 717.40 M80.88XK 736.0 M87.077 M92.30 M02.211 M25.532 737.29 754.59 733.94 726.8 M67.371 M85.572 M41.30 M08.229 M05.351 M05.512 M42 M25.229 M02.39 754.60 M42.13 M25.871 719.64 717.5 M06.272 Q66.3 M66.132 M61.239 M66.839 D08.042 M21.061 M93.952 M89.462 M96.661 M25.142 M86.8X3 M85.421 M10.131 M91.42 M89.321 M25.419 M71.451 M87.339 M70.911 524.12 M61.00 754.6 M06 M84.451G M86.329 M89.221 M80.052G M89.359 M05.712 M84.758 M86.659 M40.37 718.18 M24.451 M60.061 M47.25 M01.X32 M05.562 M84.476K 726.72 M97.02XD M14.669 M80.052P M87.236 M85.422 M86.172 M89.339 M80.822A M21.539 718.08 M41.43 M84.434 M05.672 754.7 M84.431P M85.839 M05.052 M67.829 M84.663D M75.32 M90.859 M87.363 M12.369 M84.476S M12.062 M80.019 M06.331 755.35 M12.079 M60.162 M84.532 M84.633D M99.16 M06.09 M24.312 M95.4 M90.662 K45.8 M46.05 M08.012 M11.272 M61.19 M46.24 M84.671 M84.673D M48.8X9 M79.9 M84.755S M12.39 M84.58XA 736.73 M21.511 M00.271 M21.162 M66.329 M47.14 M84.443P 719.01 M05.741 M66.331 M84.659G M67.429 M20.22 M80.831P M84.431A M60.045 M25.339 M60.242 M20.62 724.01 M19.279 M48.56XG M65.172 M48.41XG K08.26 M71.821 M06.012 M66.322 M12.09 M66.369 M86.279 M47.813 M11.271 550.11 M12.462 M54.14 M05.122 M79.645 M71.039 M89.231 M87.372 M62.449 M11.111 M05.862 M92.01 Q72.10 M47.24 M99.42 M43.22 M20.10 M06.341 M33.29 736.72 M56 M71.819 M86.452 M65.28 M00.842 M50.80 M79 755.33 550.93 M80.039K M00.152 M25.271 M66.122 718.93 M50.120 M84.68XP M89.59 737.39 M25.772 M19.141 553.00 M65.151 M05.022 M80.811A M70.89 718.83 M10.011 M10.261 M90.50 M11.061 M89.158 M19.039 M86.132 M16.7 M84.350D M41.07 524.63 M46.1 M15.0 M12.131 M89.272 M93.832 M48.31 D48.1 M11.18 M05.771 M35.6 M12.141 M24.652 M87.029 M61.171 M33.11 524.61 M25.08 |

| | | |
|---|---|---|
| **Neoplastic** | | C21.2 C05.8 D04.20 C81.30 202.64 172.5 C93.Z1 D46.C 277.84 C44.301 229 C60.0 C75.3 C44.691 D28.2 C66.9 D37.030 160.0 C91.00 D07.69 D3A.023 173.01 233.30 C34.01 C85.80 D44.4 C45.0 234.0 217 D41.8 209.3 273.3 C44.292 D30.8 227.9 194.9 D06.0 237.2 171.9 202.53 173.41 D29.22 206.2 151.5 C44.112 D38.5 D27.1 C45.7 600.1 D02.4 D07.2 C45.1 210.5 201.72 174 C16.8 155.2 C85.29 155.1 C57.12 192 D44.12 C50.122 C82.02 140 C81.77 149.8 C91.01 C25.8 201.50 237.3 C81.40 C08.9 C84.42 140.6 C83.19 152 186.0 198.8 D23.12 181 146.1 C44.622 211.1 C66.1 202.55 C78.6 200.50 C91.40 236.7 195.0 C96.2 196.0 196.6 173.92 C49 216.8 205.11 D23.21 198.7 206.82 C93.32 189.2 193 D05.01 173.9 140.1 200 200.21 C44.92 233.9 D30.4 C26.0 236.6 C84.07 C44.199 202.2 238.74 228.03 C10.2 D3A.020 C82.44 229.0 C83.72 150 235 C62.00 C33 202.31 D46 C78.7 144.9 188.6 C85.21 C50.422 218.1 C73 C06.80 C69.62 D31.00 D36.13 C76.41 172.2 201.74 187 172.9 173.82 161 D25.1 C44.82 C50.129 C48.8 200.16 175.0 D35 C25.9 C81.04 C84.44 204.10 C44.310 205.32 C49.0 I82.C29 C50.629 200.55 C49.8 D41.00 D05.82 204.81 D48.4 203.1 C81.90 C38.3 C68.8 623.7 209.34 208.91 239.5 C83.86 185 164.0 164.1 C82.46 151.8 201.0 C34.82 235.8 C84.13 198.0 203 202.06 C51.2 200.48 157.1 624.3 D42.1 C34.2 D21.4 204.90 205.91 C44.721 C69.40 164 C82.28 200.5 D01.49 C90.11 C49.A4 C43.60 D02.21 C44.219 C84.08 D30.00 C72.32 D15 C74.01 206.92 C83.13 C7A.020 C81.35 D23.61 202.42 160.1 D31.50 143.1 D31.32 C62.10 200.54 202.43 C81.99 C81.21 D36.10 D35.01 C84.A3 D10.4 C82.23 209.65 C72.21 C90.20 D05.00 C22.1 213.1 Q85.00 D24 173.52 202.11 C34.31 173.6 D13.7 D03.70 190.3 201.08 D15.0 C76.50 D12.6 D41.21 D3A.098 D26.0 184.1 208.2 C84.Z4 C84.10 225.4 C83.79 D39.9 C82.34 145.4 C69.01 C79 201.75 C84.A2 C41 143.9 C84.A9 214.4 D47.02 208.8 227.4 206.1 C15.9 C91.Z1 200.37 C17.3 C83.05 207.2 600.9 C50.222 D01 623.0 C81.20 209.13 D06.7 C09 201.73 C90 C88.8 200.30 C84.90 C40.22 C92.30 201.58 202.57 155 C50.619 D47.Z9 200.72 C82.94 170.2 C34.30 215.8 C40.00 C83.51 C92.41 201.94 204.8 173.32 173.91 E31.22 C50.311 D13.5 200.65 207.21 Z15.01 201.6 C16.6 C50.911 206.81 D31.10 D37.6 D07.60 C60.2 C82.24 C82.89 C43.39 C69.11 C82.21 D31.91 C76.51 D05.12 C83.34 C44.319 159.0 232.4 C84.61 C44.41 201.52 D03.4 C92.61 153 237.70 D21.9 D21.5 D38.2 D30.11 D18.03 209.11 233.39 D22.30 201.54 173.59 C48 209.30 C67.9 C72 C7A.00 C82.88 C16.0 C25 C18.2 C18.4 D23.72 189 C88.0 D05.02 C24.8 C68.1 C76.40 C63.12 147.0 C47.20 C82.33 C84.A5 D28.9 C91.92 201.1 205.80 C49.A2 C51.8 C78.01 600.20 D16.5 216.0 228.0 C44.91 149.9 C93.90 153.0 C46.50 200.10 199.1 D31.01 202.50 C82.81 223.81 622.1 233.6 C50.822 202.83 C44.01 C82.38 C83.91 C92.92 N60.89 C62.90 C44.712 C26.9 228.00 144.8 220 D12.5 759.5 C85.25 624.01 D03.62 D30.20 209.50 216.9 239.7 C84.Z2 C50.922 190.0 C44.399 C85.96 C85 201.07 D28.7 207.0 157.9 D3A.00 152.2 201.96 205.01 D43 C85.27 C82.30 201.28 D09.10 C91.41 C64.9 C22.4 238.76 224.5 C50.511 186 202.52 203.10 C7A.092 201.60 209.52 D41.20 C13.1 224.3 C84.04 C81.78 C90.00 D31.12 D01.3 C7B.02 209.31 157.8 C50.611 D12.1 236.99 D09.22 D17.0 202.85 C22.0 C77.9 D37.01 D06.1 D37.5 231.2 224 209.61 C84.97 C50.812 C00.0 C50.019 C84.72 173.4 222.0 C7A.022 D10.39 205 C81.32 C83.01 C83.77 227.6 C79.32 D37.2 D23.0 200.27 C85.95 621.30 C90.10 C10 200.31 C79.11 C18.8 N60.49 C89 C91.Z2 C72.9 209.02 200.67 140.5 C13.2 204 D31.40 235.5 D07.39 238 C29 C92.A2 D08 C54.0 D17.22 I82.C23 153.7 C81.05 C41.9 C00 183.3 D04.61 202.38 C06 174.4 D17.71 N40.1 D09.20 225.2 C44.500 D16.6 C84.A7 C82.95 C81.39 233.1 141 188.4 C69.42 C82.93 202.01 C83.02 C82.01 C65.1 239.4 192.8 160.2 C46 142.2 C49.3 218.2 200.2 C62.01 201.47 C84.05 C62.92 C60.1 192.9 C90.02 C82.25 202.76 600.00 C50.621 212.7 201.44 153.9 C82.61 203.11 202.97 208.01 D37 C55 C82.49 D21.10 208.12 D48.9 C22.3 C44.609 201.10 C54.1 160.9 C13 D31 C82.39 C84.01 197.7 600.0 175 C84.66 D01.2 D29.20 D35.9 D36.17 D39.2 C88.9 198 C56.9 D14.2 156 E31.21 200.68 C84.64 600.2 D37.9 239.9 C81.25 C85.16 C92.02 C44.529 D31.41 D12.3 227.1 212.8 221.0 C44.90 C50.111 C81.48 C84.92 C53.8 209.20 174.1 C00.8 202.16 237.7 205.82 C56 Q85.03 C84.45 173.70 151.4 D23.71 C83.18 C92.Z2 'C81.78 205.30 239.2 C96.29 C50.11 224.3 D16.4 C16.3 C50.829 C71.8 C92.31 200.78 C71.5 230.6 190.1 202.14 C82.16 201.00 C69.52 D12.9 173.89 D15.9 205.8 202.28 D29.21 D19 C82.57 C84.62 C43.0 C81.70 C94.01 C81.08 D02 214.8 C60.8 C83.06 C13.9 231.8 C72.59 C90.32 C50 187.8 D14.1 202.88 C51.9 202.81 C83.36 C44.509 C44.619 D04.0 D26.9 D37.8 214.3 C25.0 C46.2 182.8 C75.5 D07.30 200.38 C70.9 D15.2 188 C84.A1 C44.391 C84.02 202.26 173.81 C88.4 C91 D29.31 C94.6 C57 202.4 200.34 Q85.1 218 236.9 173.11 C18.7 C31.3 201.25 C14.8 D22.22 D35.5 C71.6 C90.21 C43.62 C49.6 236.5 188.0 C31 C22.9 C50.212 C57.7 C93.02 171.0 201.76 C00.5 C25.1 213.4 200.57 C82.64 C84.98 D35.3 D49.512 C67.1 D04.10 C83.14 C75.9 200.23 200.26 161.0 C78.1 200.7 C82.85 D31.22 C83.75 D21.11 200.83 200.43 C57.4 D10.7 146.5 D49 C44.510 C57.11 238.0 C50.021 196.3 D46.2 D16.31 C78.39 569.0 D39 161.2 C25.3 C30.0 C01 D31.31 D36.14 201.46 C94.81 140.4 150.3 C91.A2 D16.02 D13.2 C47.10 C26.1 160.8 D23.10 173.8 D37.1 C44.701 C11 D25 C63.8 C94.00 D17.24 209.74 C02.9 C91.61 C28 205.02 D21.0 171.6 C17.9 210.9 174.0 D21 C67.4 C50.821 C49.5 173.7 C92.11 209.73 C84.67 C74 C7A.095 219.8 161.8 C06.9 202.75 200.80 224.9 C84.73 C88.2 D17.79 201.21 208.90 152.1 D03 147.8 186.9 C50.622 198.6 229.9 C49.A5 D22.62 203.80 C40.32 C63.9 C81.33 C96 209.6 C87 202.86 C43.30 207.01 D35.1 164.8 238.73 201.62 C96.21 D04.71 171.4 201.41 141.0 C91.60 I82.C21 206.80 189.8 208.21 C82.03 D03.52 D18.09 C25.7 162.0 C38.0 202.72 C82.41 C82.62 173.39 C72.41 202.00 C03.9 229.8 173.40 224.4 213.3 216.3 209.14 C83.8 C16.4 C75.4 202.05 C44.191 213 200.13 D41.11 C69.50 192.1 200.75 D13.1 D48.3 C00.2 C80.0 C81.09 C83.89 195.8 C93.Z0 N40.0 C30 C92.A0 C00.1 C34.10 200.1 174.2 173.71 158.0 D07.61 C84.91 156.1 C84.12 C91.42 D00.07 221 D35.6 225.9 624.6 233.31 D19.1 C91.A0 C05.0 192.3 211.4 C82.45 160.3 622.10 211.5 C7A.011 D10.0 204.0 C81.42 201.26 D23.30 C02 201.91 C40.90 C78.02 C43.21 C44.222 C91.62 239.3 237.71 239.6 230.8 202.8 232.0 C44.89 C81.12 C7A.021 C96.Z 200.01 202.51 N89.3 173.09 D27.9 156.9 C94.82 148.0 197.4 D36.7 202.02 C03.0 189.9 201.53 205.0 D00.02 273.1 199 200.85 E71.448 C50.121 202.20 D48.2 C62.02 149.1 C83.00 C50.321 173.99 D44.2 173.1 C22.7 D44.6 208 C82.56 157.0 170.7 D24.1 D14 C71.4 600.90 C85.98 |
| **Ophthalmological** | | H34.8332 H50.812 H10.439 H35.111 H53.043 364.5 H50.32 H44.812 H01.112 H05.243 H26.493 H10.221 377.62 362.89 H49.43 H18.811 364.71 375.13 365.59 365.89 370.3 H40.51X H44.21 370.24 375.14 H25.23 H17.00 H47.631 H05.342 H21.323 H50.811 H59.022 364.8 H21.229 H20.022 H02.134 H15.022 H46.02 H47.211 H31.111 H59.322 H35.3133 H11.032 H11.432 H16.072 H40.10x2 371.73 370.03 H11.10 362.31 H02.734 H44.753 369.3 H16.142 H40.63X3 H33.051 H18.739 369.05 360.8 H31.092 H52.532 H18.892 H17.03 H04.223 H01.133 371.20 H27.112 H44.621 H20.813 H53.033 367.20 369.08 372.9 H02.871 H40.1310 H40.51X0 H15.091 H04.302 H01.023 H31.303 H26.119 H40.1114 H54.2X1 368.31 H15.813 H35.463 H05.812 H46.011 372.4 364.53 H10.011 372.4 347.333 H01.019 H21.342 H02.30 H20.821 H40.53X3 H02.104 H02.839 H05.031 H59.343 H16.061 H16.071 H34.02 377.10 371.45 371.60 H18.623 369.75 H02.209 H33.102 369.60 H50.50 371.51 H40.032 H40.2214 H33.22 H05.53 H05.813 H44.649 H54.1224 H02.439 369.66 H40.049 H40.51X1 H40.10x H02.054 H34.00 H50.15 H18.55 H15.821 H44.511 H49.813 H40.53X0 372.56 H02.136 H21.1X9 H53.023 H04.123 H02.89 H18.422 H31.123 H02.036 H20.021 374.54 368 H26.412 H40.60X3 H46.9 368.12 H57.13 H25.811 360.40 H40.62X4 371.58 369.01 H10.413 H16.309 H26.042 H44.431 H49.02 H53.031 362.6 H11.152 H35.3114 H26.3 H44.131 H35.3114 H53.133 366.19 H40.1111 H08 H44.2C2 H53.433 H05.023 H50.69 H04.541 360.34 H04.212 H05.263 H11.133 H31.411 H18.232 H30.111 371.54 H04.023 H40.61X H35.172 H02.203 362.24 H33.032 H40.32X4 H44.712 377.52 H21.319 H20.22 H02.409 H10.12 H15.812 366.10 H18.9 H52.201 H50.022 H11.423 H34.813 H54.60 H57.051 H18.062 H40.63X H44.2A3 H44.539 361.10 H27.133 368.16 H40.233 H30.133 H43.813 H02.109 H18.712 H57.8 371 H44.131 H35.3114 H53.133 H40.42X3 H53.133 366.19 H01.111 H08 H10.223 H40.139 H18.312 H35.322 H49.01 371.0 H02.735 374.87 H30.012 H02.123 H44.2C3 H13 H44.2E2 366.1 369.62 H27.119 H31.21 H47.323 H53.422 364.7 H16.051 371.56 H35.079 H40.1312 H53.039 H15.052 H40.1290 H30.009 H40.20x2 H01.002 H40.041 365.21 H34.819 375.54 H02.043 H44.449 H54 H18.713 362.85 364.59 H44.629 H46.12 H47.032 H40.159 364.63 H02.006 H40.822 H21.219 H44.022 H35.3290 368.41 H05.339 377.11 H05.269 H11.063 H54.1152 360.81 H05.52 366.04 H52.13 H53.429 362.77 369.61 369.10 H01.013 H02.879 367.9 365.04 H35.413 H11.242 H30.91 H11.019 H02.239 H30.029 H35.70 H40.43X1 H59.312 371.11 368.2 H04.149 H18.022 376.42 H40.41X2 368.8 H15.819 H05.262 H00.10 H16.212 H59.88 H21.532 H01.125 H31.32 H02.714 H05.423 H21.521 H40.053 H54.0X4 362.70 H02.61 H10.231 H51.22 H35.013 H34.8131 369.65 H52.12 370.40 H21.352 H35.3232 H15.123 H02.832 369.24 H57.11 H59.42 369.02 H50.611 H44.441 H31.009 H34.8392 H54.414A H15.841 376.5 374.13 H16.121 379.14 H35.359 H44.811 H54.1141 H47.611 376.21 H04.432 H05.219 H40.123 H21.262 H34.832 H44.702 H00.13 H40.221 H43.392 H40.1412 H01.129 H44.652 H02.863 H40.50X H00.034 H30.101 H16.212 H59.88 H21.532 H01.125 H31.33 H22 H33.029 H02.721 362.10 H20.041 362.3 360.44 H40.52X4 H02.013 H21.40 H07 H30.129 372.8 375.15 H21.232 H01.016 H31.029 365.05 H40.1491 H53.59 H35.3121 365.6 H02.32 367.2 H11.212 H44.622 H05.021 H50.312 H16.032 H02.516 362.14 H40.1123 H43.312 372.5 H44.732 H04.551 372.42 H44.512 H33.322 379.25 H05.231 H47 H50.30 H21.301 H21.343 H59.211 H40.42X2 H16.323 366.20 H40.62X2 H10.411 H20.042 H04.013 H15.102 H35.411 H16.019 367.5 371.3 H40.1193 H59.319 H18.812 H35.019 H47.013 H02.432 H40.1291 H53.142 H21.333 H54.512 H10.021 369.04 362.16 369.72 H02.875 H30.021 H43.01 365.72 H59.813 368.55 362.74 H31.311 H30.103 H47.219 363.56 H02.833 H40.1194 H00.035 H04.211 H47.393 371.24 379.32 H02.025 H02.204 H15.832 376.41 H15.833 376.1 H16.332 H40.51X4 H43 H40.52X1 H11.422 H40.1423 H40.60X H52.203 369.0 374.89 H02.849 374.12 H40.50X0 H54.8 H40.043 H44.611 365.60 A18.59 362 368.34 H50.011 H59.021 369.71 367.0 362.15 H11.32 H18.319 H11.131 H53.483 H02.59 H35.3291 H40.30X2 H44.529 H35.319 H40.40X2 H11.052 374.5 H16.391 H21.263 366.42 366.18 H16.129 H40.1221 H18.419 362.29 368.47 H21 365.23 376.40 H18.819 H21.321 H40.63X1 H30.21 H50.00 H33.052 379.21 H02.514 H40.1121 H47.42 H53.411 H15.019 H05 H54.40 372.55 H40.1120 H35.113 H05.011 H59.222 370.44 H33.192 H59.333 362.61 362.17 H02.522 H02.119 H10.812 H59.323 H54.0X45 H18.039 H43.02 H35.722 366.09 H21.313 H04.523 H21.551 H05.033 H54.62 H54.42A H58 H04.332 H10.502 H34.829 365.31 H42 367.1 H30.013 375.57 H21.303 H47.20 H52.864 H20.032 H26.062 377.3 H52.223 H31.002 H40.2232 H44.623 H44.821 H27.03 H04.012 H16.013 H34.833 372.62 H02.055 H02.114 H30.121 H27 377.33 369.9 H16.223 H16.213 H18.829 363.72 H35.723 H25.041 H44.2E3 H44.012 H34.10 H16.221 H40.069 H02 363.55 H04.202 H16.131 369.06 H31.9 H35.3222 H34.822 H35.3293 H47.512 H44.413 H20.21 H16.331 368.62 H35.371 H40.41X3 H21.312 H33.331 H16.012 H33.331 H16.012 H33.331 H80.893 H26.499 H31.429 H47.12 377.24 H35.3120 H40.1130 H18.719 371.9 364.52 H40.1322 372.34 H11.411 H16.073 H10.232 H10.212 376.43 H35.81 H35.712 H01.025 H11.069 362.65 379.13 H01.123 H35.373 H16.122 H40.2220 H43.811 H31.409 362.1 H34.213 H02.815 H01.029 H40.539 H05.412 H52.521 H21.242 H02.714 366 H34.8191 H34.8391 H35.53 H21.89 H01.001 377.41 H49.33 H31.321 H44.023 H16.052 369.17 370.21 H16.249 H55.09 H26.051 H18.061 H21.223 H25.11 H30.032 H59.362 H53.10 H04.321 H40.1234 H20.12 H11.111 H04.119 H40.1113 H04.031 H21.213 375.52 H26.491 H21.1X1 H40.2292 H26.011 H26.033 H43.20 H15.041 368.9 362.50 H35.349 H15.89 H52.229 H51.8 364.64 H02.829 H00.039 H18.322 H40.42X1 H40.812 H50.42 H02.034 H05.029 H21.561 372.42 H40.40X3 H40.832 H40.52X3 H05.10 H54.413 374.51 H05.012 H40.1223 H11.829 365.32 377.03 H40.1432 368.32 371.05 H11.431 H01.124 368.30 H34.233 366.4 379 H10.819 H53.129 H02.731 H02.205 H04.429 H49.32 371.02 H53.15 H40.1402 H40.234 H40.022 H33.011 H35.732 H44.391 366.23 H04.039 H50.17 H16.139 368.52 H44.793 H11.059 H18.449 H18.231 H05.823 H18.032 H35.352 H26.049 H04.129 H10.30 H02.811 H11.113 H02.212 H11.439 366.2 H35.3292 H35.441 H40.10X0 H35.023 376.51 H31.103 H43.391 H02.016 H35.3110 H05.253 H05.821 H40.219 H05.409 H53.021 H21.329 H04.439 365.11 H17 H49.11 368.3 372.54 361.3 H54.52A H25.093 H02.145 H44.522 H04.132 H11.022 H04.552 H16.041 H44.722 H54.415A H38 H40.1213 H01.121 363.9 H02.403 H40.40x0 H35.461 365.35 H18.793 376.8 H35.292 H35.441 H40.40X0 H35.023 H48 H18.52 H34.8390 H40.51X3 H04.339 H47.11 H01.022 H05.811 H40.141 H04.301 H31.422 H15.823 H21.503 375.69 H40.009 H01.146 H57.052 H30.139 H44.601 H26.411 H27.132 H18.899 366.41 H54.61 H20.819 H54.42A4 370.34 H33.129 H21.1X3 H02.019 H26.112 H44.323 374.41 H05.221 363.33 366.43 H18.421 H52.11 376.47 H53.131 H44.602 H30.031 364.60 H10.10 H11.229 H00.021 H16.063 H02.133 369.11 H16.011 H52.32 H18.20 H31.403 H52.533 H04.612 H31.23 H43.311 H35.459 372.5 H16.233 H35.713 H11.239 H16.412 H40.2293 H33.121 364.51 H35.012 H53.439 H40.213 368.40 H18.229 H25.011 H01.021 H17.13 H30.143 H02.812 H44.2D1 H35.341 H04.133 H18.799 H35.3192 H59.329 H26.213 H18.54 376.89 H01.011 H31.102 H21.233 H16.422 H53.19 H04.522 H31.29 H47.233 H53.421 368.60 H49.31 H31.412 |
| **Osteoarthrosis** | | M17.5 M19.079 M19.239 715.24 715.25 I20.8 715.33 M17.9 715.00 715.35 M17.10 I25.3 I21.11 M15.1 M19.249 M19.029 M19.279 M15.3 M19.91 I25.42 715.96 715.21 715.94 715.34 715.14 715.11 M19.039 715.32 715.23 I25.10 715.10 715.92 I25.811 715.12 I25.2 I21.29 M15.0 I21.4 715.27 715.18 M18.9 715.04 I25.83 715.28 M15.8 M16.9 715.89 I25.9 715.22 715.30 M19.93 I24.1 I25.41 M19.90 715.31 M15.9 715.80 M19.019 715.09 715.13 I24.0 M16.10 715.36 715.20 M16.7 715.93 M19.219 I21.19 I20.0 715.37 I20.1 715.15 715.17 I24.8 715.98 I25.89 715.95 I25.810 715.97 715.16 I21.3 M19.049 I25.82 715.90 715.91 715.38 M19.229 715.26 I21.09 I25.812 |

| | |
|---|---|
| **Oth-Joint-disord** | 719.35 719.05 719.13 M25.40 M25.459 719.07 719.86 I20.8 719.33 M25.18 719.91 I25.3 I21.11 719.54 719.58 719.62 719.89 M25.073 719.29 719.38 719.27 719.25 719.92 M25.619 I25.42 M25.639 719.30 719.53 719.85 M25.48 719.90 M25.559 719.14 719.47 719.42 M25.439 719.34 719.88 I25.10 719.22 719.36 719.94 719.56 719.96 I25.811 719.81 M25.00 I25.2 M25.059 M25.08 719.99 I21.29 719.82 719.03 719.68 I21.4 719.11 719.09 M25.039 719.01 719.28 719.15 M25.019 M25.519 M25.60 719.60 M25.069 719.21 M25.029 719.31 719.66 M25.449 719.43 I25.83 719.08 M25.9 M25.879 M25.849 M25.649 719.93 M25.10 719.02 719.80 719.98 I25.9 719.04 719.32 I24.1 I25.41 719.18 719.55 M25.629 719.7 719.46 719.52 I24.0 M25.659 M25.669 719.19 M25.049 719.67 M25.579 M25.859 719.37 I21.19 I20.0 719.61 M25.80 719.24 I20.1 719.57 719.23 M25.429 M25.839 719.00 M25.469 719.84 719.97 719.59 719.39 719.45 719.65 I24.8 719.10 719.50 719.63 M25.473 M25.729 M25.70 719.26 M25.539 I25.89 719.16 719.44 719.41 719.49 719.83 M25.869 I25.810 719.87 M25.829 M25.50 719.17 719.20 I21.3 719.95 719.48 719.51 719.64 M25.673 I25.82 M25.119 719.12 719.69 M25.419 719.40 719.06 I21.09 I25.812 M25.529 M25.569 |
| **Oth-Urinary** | N39.41 599.71 N39.45 N39.42 599.4 I24.8 I25.10 599.69 599.3 I25.89 I20.8 599.5 N39.44 I25.9 599.60 I25.811 N39.8 I24.1 I25.2 N39.46 I25.41 I25.810 I25.3 I21.11 599.70 N39.43 I21.29 I21.3 N39.3 599.9 I24.0 599.0 I25.82 I21.4 599.2 I25.42 599.80 599.83 599.82 599.72 599.84 I21.19 I20.0 N39.498 599.1 I20.1 I21.09 I25.812 599.89 N39.0 599.81 I25.83 |
| **Otic** | H66.006 H60.23 H93.249 H62.41 H91.02 388.71 H61.391 H80.20 H93.219 H93.212 H93.099 381.89 H95.811 H72.10 H73.91 H74.09 H66.21 H70.229 H90.A12 H95.123 384.00 H60.312 H61.302 H95.193 H65.199 H61.011 389.7 H70.009 380.39 H71.31 H69.91 H60.543 H94.01 H66.3X3 H80.91 H75.00 H65.191 H73.822 H68.019 385.8 H65.04 388.0 H68.112 H68.131 H92.20 H60.61 H95.54 H90.2 H74.392 H60.532 H70.211 H65.20 389.17 H69.90 H95.89 H90.3 388.70 H71.33 H74.41 H65.413 H83.12 H90.6 H67.3 H60.8X2 H72.93 H61.323 H80.23 H60.323 H80.92 H72.829 H90.A21 H72.01 H92.21 H93.13 H60.531 H61.103 H74.321 H72.11 H80.01 H70.092 380.5 H95.119 H72.813 381.9 H81.313 388.02 H81.43 H66.10 H70.811 H61.23 386.40 H83.90 H94.03 388.72 H94.80 H90.12 H74.40 386.32 H83.8X3 H73.811 H75.02 H95.121 H60.539 H73.093 385.13 H66.3X2 H90.A32 H61.003 H60.02 H73.012 H91 380.21 380.30 H95.03 H95.53 H83 386.31 H73.001 H95.111 H60.502 H88 H66.005 H65.01 H66.42 H81.12 H60.8X3 H60.42 388.2 389.02 H61.392 380.4 H61.031 H60.43 H91.8X9 H70.099 381.8 H66.009 380.03 H83.3X9 H93.293 H68.012 H62.8X2 H74.42 H72.12 H70.202 380.0 H74.8X3 H60.60 H61.113 H61.819 H93.243 386.48 H65.31 H73.003 380.00 H82.3 H95.88 389.03 H65.06 H79 H60.03 H81.41 H61.393 388.60 H92.13 388.00 H83.3X1 H60.10 389.0 H62.43 H83.19 H95.133 380.50 H61.019 H81.01 H93.091 H73 384.1 H66.23 H71.10 H74.391 H80.03 H72.02 H60.599 H66.91 380.51 380.52 381.62 H60.592 H71.32 H70.011 385.0 H73.099 H66.12 H73.091 H68.023 H93.232 H93.8X2 H77 388.11 H72.823 H60.00 H70.093 H65.93 H81.93 H67.2 H93.3X9 389.16 386.5 H60.559 389.08 H73.019 H61.013 H60.512 385.01 H95.113 H70.209 H81.49 H65.21 H74.323 H73.813 H70.012 389.14 H69 H62.42 H61.001 H90.A22 386.43 385.23 H61.102 H95.21 H68.129 H83.01 H61.129 H65.196 H72.90 H65.32 H83.2X2 384.8 H73.22 H95.139 H65.90 H74 384 H72.00 H91.20 H83.8X2 H71.23 H89 H93.3X3 H62.8X9 H73.829 H66.11 H71.00 389.11 H65.03 H74.399 H69.03 385.11 384.23 H60.511 H90.72 380.02 H60.391 H81.8X2 386.30 H91.09 H95 H66.003 H75.83 384.25 H83.02 H65.115 H73.10 H65.197 H60.541 H61.811 385.89 H60.331 H81.92 H70.219 385.24 380.9 H61.399 H68.022 384.21 H74.92 380.81 H60.40 H66.93 H92.23 H65.492 H69.02 H70.893 380.8 384.01 H73.821 389.05 H73.93 H91.01 H65.05 H74.393 387.1 H72.13 H81.21 H81.391 H73.013 H70.11 H60.399 H61.123 H69.80 H74.312 H61.93 H74.322 H83.11 H91.91 H75 383.32 H91.93 H74.311 H73.11 381.5 H93.241 H80.83 H90.5 H92.10 385.2 H61.039 380.89 H95.42 H68.133 H80.11 386.34 384.20 H81.11 384.9 H73.819 H93.A9 H80.90 H65.116 H65.193 H73.92 H65.113 H60.542 H70.813 385.21 H90.71 H81.23 386.58 385.12 380.31 381.6 H70.92 H66.22 H65.114 H80.93 H92.03 H65.419 H61.111 H93.223 H93.A2 H94.82 H60.521 H66.012 381.81 380 H93.12 H95.129 H93 H90.A11 H74.02 H81.22 385.35 H71.90 H93.239 H69.01 H95.02 386.56 387.8 389.06 H66.014 H70 H73.011 H60.333 384.82 H60.20 385.00 H60.392 385.09 H61.032 H74.329 H83.09 388.31 H74.22 H74.03 H93.019 H71.92 388.3 H70.12 H61.311 H93.221 H93.013 389.01 H61.023 H66.019 381.63 H60.529 H71.91 H95.813 H60.549 H75.82 H81.399 H61.119 H91.8X3 H65.22 H93.211 H70.212 H70.213 H60.92 H71.12 H61.92 H73.21 H61.303 H80.21 H61.303 H73.091 H66.20 H71.11 H72.91 H90.0 H61.121 H70.222 H95.819 H61.301 H70.891 H73.891 H73.009 H65.07 H66.007 H61.029 H68.113 H71.02 H82.9 H61.813 H70.001 H81.393 H93.8X9 H61.009 H61.109 H94.02 H61.319 386.41 H93.012 H60 H60.8X1 H80.12 H81.312 385.33 H70.812 H74.12 H70.91 H95.132 H93.222 381.52 389.8 H80.10 H66.001 H81.8X2 386.30 H93.299 H80.80 H65.195 H68.122 388.01 H95.22 H66.40 H74.13 H95.191 H81.8X9 H60.591 H61.892 H72.2X1 H74.313 H61.329 H65.30 H61.193 H81.13 380.32 387.2 H68.102 H60.523 H66 H74.43 H60.501 H83.13 H93.90 387 H91.3 H94.81 H95.52 H66.92 H80.02 H70.13 H73.893 H73.892 H65.00 H90.42 H60.62 H93.92 H60.41 H65.411 H91.11 H91.92 H71.30 389.9 H60.313 H62.8X3 H70.091 H83.2X3 389.00 H60.63 H74.20 389.12 385.82 H65.192 H93.092 H72 H68.109 H65.92 H73.90 H80.22 H69.00 H60.322 H66.004 H90 H68.132 H60.503 H71.01 H68.111 H91.13 H74.8X1 H64 388.6 H95.192 H72.821 H81.8X3 388.7 H61.012 H83.3X2 H74.01 H95.131 H66.3X1 H93.11 H83.8X1 386.50 H75.81 386.55 H61.112 H65.33 H60.332 H72.2X3 H81.02 H72.03 388.9 H60.11 380.01 H61.20 389.22 H62.8X1 H95.41 387.9 H73.812 388.8 H60.522 389.2 388.69 H91.8X2 H95.32 H95.812 H92.12 H93.291 H70.10 H70.201 H71.22 387.0 H71 385 H72.2X9 H61.321 H65.194 H66.002 H61.91 H92.09 H61.313 H83.2X1 H66.43 H86 H83.8X9 H91.21 H70.002 389.13 H60.93 H67.1 H69.81 H81.42 H90.11 H71.21 H78 H80.13 H84 388.30 H61.812 H61.021 381.7 H74.8X9 H60.311 H74.90 H90.A31 H60.12 H60.8X9 H69.83 H65.493 H60.551 384.09 H61.893 385.03 H93.3X1 H60.513 H66.123 H72.822 H68.002 H63 H60.393 H91.8X1 H93.011 H61.891 H66.013 H70.899 H60.329 H61.309 H80.82 H62.40 H60.13 H66.016 H75.01 H91.23 H81.09 380.3 386.33 H80 H66.011 381.61 385.3 H93.093 H61.899 H81.91 H93.19 H81.319 H74.21 388 H92.01 H68.021 H83.92 H72.819 H92.02 H71.13 H61.192 H60.321 H68.009 H61 H74.319 H93.8X1 H68.119 H81.8X1 H68.001 H70.90 H83.03 385 H61.033 H65.491 386.54 H93.3X2 H70.003 H75.80 389.20 H93.A1 389 H68.013 H81.90 H60.01 H75.03 H72.811 H66.3X9 H70.203 H81 381.60 384.22 H95.199 H81.311 H92 389.04 H65.117 H68.101 381.50 H70.93 H73.823 H80.00 H60.519 H83.2X9 H61.122 H74.8X2 H60.509 386.52 385.9 H61.002 H82 385.02 388.32 H95.01 385.10 386.51 H61.191 H60.593 H73.13 H82.1 H95.31 H61.312 389.21 H74.23 H93.A3 H70.819 388.12 386.4 H60.22 385.19 H65.499 H69.82 H71.03 H74.93 H68.029 H70.013 H91.12 H70.019 H74.91 H70.223 H83.93 H81.392 386.53 H60.553 H61.101 H93.25 H60.91 H67 385.1 H68.103 H71.20 H65.23 388.10 H69.93 H74.19 H66.017 H73.20 H68 H90.41 H66.41 H91.03 388.1 H60.552 H95.31 H82.2 386.35 H61.322 H60.90 H81.03 H68.139 385.30 H65.91 389.18 H93.292 H95.112 H94 381 H93.233 386.42 H94.00 389.10 H72.92 H70.221 H72.2X2 H73.899 H93.242 H66.015 385.22 H92.22 384.0 H73.092 H93.91 H76 H69.92 H61.022 384.2 H91.90 381.51 H73.12 385.31 388.61 386.8 H93.213 G96.0 H66.90 H61.22 H68.121 |
| **PNS** | G56.91 G73.3 G81.92 G62.82 355 359.21 P11.5 G81.11 352.5 G57.53 352.4 G82.22 G83.14 G57.00 354.3 G80.0 354.4 355.3 354.2 353.3 G83.10 G50.0 767.6 R26.9 G56.43 G79 G71.19 G60.9 G83.13 350.2 G81.02 G71.13 G80.90 G52.9 G57.62 G59 R26.1 R26.0 356.2 G51.8 G57.92 P14.8 G57.21 G54.4 G72.2 357 G83.30 359.22 358.2 781.94 351.8 G71.11 G54.7 G72.9 G57.22 G83.89 359.81 G56.42 G60.2 G72.41 G82.50 G80 G56.02 359.71 G82.20 359.79 G83.32 G57.10 353.0 781.1 357.7 G57.80 G81.01 359.4 G57.43 354.8 G57.40 G56.00 G59 R26.1 G80.2 G81 G57.61 G56.80 355.2 353.4 G52.2 354.5 G81.14 358.9 G55 G72.89 G83.9 350.9 G57.83 G68 G58 358.8 G82.53 G57.82 G66 781.0 G69 355.0 351.0 R27.8 G56.83 R25.8 352.6 353.1 G70.89 G56.92 G83.84 G80.4 R29.5 352.1 G70.00 353.6 G61 G56.32 R41.4 R26.2 781.3 G83.4 R29.3 G51.0 G83.23 359.89 G82.52 G54.6 G56 G56.03 G83.20 G56.90 353.5 G60.1 G73 P14.3 G50.8 G60.0 350 358.01 G51.4 G57.02 G50 G61.82 G62 357.9 R43.0 359.3 G57.71 353 352.0 R25.0 G60.8 356.8 355.1 357.89 G54.8 G80.8 G57.90 G71.9 G57.12 G80.3 357.1 781.5 G65 R29.1 R25.1 G61.81 356.4 357.0 359.29 781.2 G54.1 G70 P14.9 G53 G56.33 G65.0 356.3 G57.73 G56.23 G83.5 353.8 G57.93 781.99 R25.2 G57.13 352.9 G83.24 G54.5 G57.72 359.9 G61.89 355.71 G82.54 355.8 G56.11 G71.8 P14.0 G51.1 351.9 359.24 355.6 G56.93 355.4 G57.50 357.2 G70.01 G76 G73.1 G57.20 359.6 G83.11 351.1 781.8 G56.21 G82 G50.9 G83 R25 G57.63 G65.2 G57.62 G54.0 R25.3 G57.31 G83.82 357.3 355.79 G56.31 G81.00 352.2 G57.91 P11.4 G70.81 G78 R27.9 G62.9 G71 G56.20 358.1 G72.0 G74 P14 767.5 P14.1 G72.3 G83.31 357.81 359 E13.42 G57.70 R27 G57.33 G70.1 G80.1 359.0 356.9 G54 G81.10 P11.3 357.5 G71.12 356.1 781 G51.9 G52.0 350.1 G58.8 G77 355.5 G51 358.00 G81.94 358 G81.04 G56.22 G83.21 359.23 781.6 G83.34 G67 781.93 G57.41 R29.891 G57.11 G54.3 R29.810 R29.890 359.1 G57.42 352.3 355.9 G52.8 G57.01 353.2 767.4 G57.30 G62.1 P14.2 G54.9 G54.2 G58.7 G52.3 R26.89 G56.12 781.91 G52.7 G70.80 781.7 G70.9 G57.23 356.0 G83.33 R27.0 G52 G61.0 G50.1 G62.2 354.0 359.5 G51.2 R29.0 G60 354.1 356 G62.81 357.6 G83.22 G71.2 354 G61.1 G56.30 G82.21 G73.7 G81.13 G56.41 G56.01 G81.12 G51.3 G83.0 G83.83 R25.9 781.4 G57.51 R29.818 781.92 357.4 G71.3 G57.52 G62.0 R68.3 G57 G63 G83.81 G56.13 G52.1 G57.03 G56.81 G72.81 G81.03 352 G64 353.9 R26 G61.9 G71.0 G71.14 G72 G83.12 G82.51 G72.49 G60.3 354.9 R26.81 350.8 351 G57.60 G57.81 G80.9 G58.9 G75 357.82 G81.91 G56.82 G58.0 G70.2 G56.10 |

## Psychiatric

F12.90 F91.8 F12.229 302.3 307.81 F41.3 F32.9 F10.29 F52.1 300.14 313.22 R44.3 F11.21 F13.221 F63.89 F60.81 R40.2440 R48.0 F16.24 F10.230 F15.251 F51.02
F35 F40.298 292.84 F42.2 R40.2254 301.4 308.4 F20.2 F14.259 F40.231 F31.4 F40.241 F17.293 296.35 F18.27 F81.89 301.8 312.10 F19.988 F14.94 F42.8 F64.9
292.12 F98.5 295.1 F80.0 F01 F13.90 F17.299 295.45 F15.250 F10.288 312.2 F13.150 F18.259 F31.0 301.9 F84.3 296.61 F13.939 302.73 293.9 F18.19 F13.232 301.12
F41 F10.150 F31.76 F45.21 R40.2240 296.81 F18.988 F19.99 F13.10 F48.1 F11.250 312.35 F40.212 F10.94 308.1 290.1 R47.1 299.10 F13.288 298.8 R40.2132 295.70
F13.99 F11.11 F30.10 F34.9 300.22 F11.14 R40.2231 315.9 307.4 F15.280 R40.2110 293.83 315.4 F9.181 300.01 F99 F80.82 R47.01 F15 F47 296.24 F96 R41.4
F12.120 F40.290 307.49 F55.4 R41.0 301.13 295.00 F14.10 307.44 F30.4 F51.4 R37 F17.218 315.31 313.2 F10.920 R40.2333 R45.3 F44.7 F25.8 294.0 307.3 R45.2
F20.3 F19.920 F20 295.43 F15.951 F10.26 316. 291.1 296.90 F16.150 295.63 292.81 296.44 315.3 302.89 R46.2 295.61 301.89 R43.0 294.9 315.1 312.30 F83 F15.259
F42.3 299 R49.22 F11.90 F15.20 F10.951 F91 F12.222 R45.1 R40.2213 F20.5 F15.11 F19.20 F19.129 302.51 F14.14 298 F31.72 295.75 F14.182 F31.13 295.25 295.20
291.82 F94.1 300.00 R40.2212 294.11 F07.0 F18.929 294.10 R46.1 F13.14 R40.214 306.5 F03 295.34 F16.920 F31.73 295.51 291.9 295.03 F50.2 F18.921 F20.0
291.89 F74 312.01 F30.12 F19.29 296.56 F32.1 F18.24 295.8 F90.0 R48.1 F13.96 R45.82 296.40 R40.2211 R40.2444 295.92 296.36 293 F45.1 F40.228 F14.980 F15.23
315.5 F10.188 F16.188 F65.9 F13.19 300.20 F41.0 309.2 F10.280 F52.6 F64.2 F18.959 F19.10 F26 F16.20 F19.21 F19.14 F65.81 F19.97 F80.2 F13.932 F10.24 F14.21
F51.13 307.40 F18.11 F02 307.0 296.4 F43.23 290.41 296 R40.2114 R40.2331 298.2 R40.234 F52.32 F31.60 F45.42 F19.980 R46.6 300.9 F06.31 F90.8 F50.02
F17.208 F13.151 F13.21 F19.932 F43.29 F45.22 F19.230 R40.2244 F30.11 R40.1 F16.283 F93.0 F43 313.89 F11.10 F14.121 F18.980 F31.32 F15.122 F63 298.4 F98.3
313. F64.1 296.99 F80.1 F09 F19.951 F63.0 F18.29 F19.250 F43.9 F10.11 R40.224 F43.25 F60.5 296.80 297.9 R44.9 F98.1 R40.2413 F31.11 F19.90 F10.20 306.59
312.3 R40.2121 291.8 295.73 F11.94 F14.922 F33.1 F25.1 F14.180 F49.8 R40.232 296.33 310.0 R40.225 302.0 F15.920 301.2 297.2 F06.4 F84.8 296.15 F17.211 F32
F17.201 F19.96 301.7 R40.2361 F19.122 F65.1 F67 R40.2342 R40.2220 296.10 308.0 300.11 F18.229 F10.231 300.29 F15.982 306.4 F18.120 310.1 302.5 292 F80.4
295.35 301 F42 F46 F17.223 F15.129 F31.77 R40.2241 F40.8 R40.2332 F59 F18.221 301.20 290.13 301.59 F12.920 306.1 F18.129 F45 F43.22 R44.2 F15.180 295.71
F18.150 F10 F18.251 R41.842 F98.29 R40.2142 R43.8 R40.2330 F13.282 R40.2314 F16.129 F19.121 R40.2341 313.23 302.1 312.23 F11.988 R45.0 F58 302.8 307.23
294.20 310.8 F10.96 295.2 306.0 F10.982 R40.231 295.22 298.1 F62 295.9 F14.288 F31.5 R40.233 F64.8 F20.81 F51.11 F06.0 F13.931 F07.81 F13.988 F77 F14.988
R40.221 F15.19 F68.8 F11.281 R40.2344 302.9 F16.122 296.03 312.1 295.11 F43.21 F84 F32.5 F13.94 F63.9 301.6 F14.19 306.2 F18.151 F16.288 F14.251 F15.221
F11 F45.0 302.76 F19 F10.988 F51.01 307.47 296.66 F12.259 R40.2424 R49.1 F90 297.8 F10.239 F51.05 300.09 315.8 291.2 307.6 F63.2 F14.181 295.01 302.81 R48
F16.99 R40.2421 F12.220 R40.2442 296.13 F11.951 300.1 F31.30 F90.9 295.83 300.33 300.6 290.4 F84.5 306.3 F33.40 295.55 F93.9 F55.0 F11.981 R46.3 297.0 F81.81 F28
F68.10 F12.950 R40.2141 290.8 295.0 F14.221 F16.251 R40.2113 314.8 295.33 F64 295.15 F48 F11.950 293.82 290.42 297 R40.2111 F84.0 R45.4 292.8 F52.31 R43.1
302.84 F14.982 R40.2251 F06.34 302.79 F19.939 F16.14 F33.0 F15.282 R40.2360 F19.259 F13.259 296.01 F94.9 R40.2232 F31.74 F15.24 302.4 F43.11 310.9 F42.4
R41.9 302.7 R47.02 F90.2 293.89 F11.922 F12.980 R46.7 301.21 296.05 F12.221 F40.240 F44.6 F16.151 F93.8 F17.203 F15.281 F68 F40.210 F13.929 F19.921 F11.20
307.48 F13.920 F11.129 296.3 R40.2441 R40.211 F16.180 309.24 F18.90 F11.121 F19.930 F01.51 F19.959 F19.17 F87 F19.11 F40.00 315.2 F76 R43 F17.291 312.
F05 295.91 F15.14 F11.182 F98.0 F55.3 306.51 R48.8 F13.120 295.41 309.23 296.06 F11.188 F19.282 F12.288 F44.89 F18.14 300.81 296.42 R40.2222 F45.20
R40.2242 F13.982 R40.2311 F15.21 F51.3 R41.2 314.2 F13.97 300.15 F11.29 301.83 F19.232 F50.81 R45.84 R40.2124 F40.01 F91.1 F10.959 F12.959 R40.222 300.82
F10.221 F42.9 R45.5 290.3 F13.27 290.2 292.89 F51.8 F15.959 296.65 F13.959 306.52 F33 F19.16 F11.222 F02.80 F52 F14.29 F31.78 294.8 315. F31.70 315.39
F15.921 F33.8 293.1 F16.90 R40.2112 R40.2323 295.54 302.2 308. 296.50 F07 315.00 296.46 F14.122 312.11 295.81 F18.159 296.16 296.00 F10.97 F07.89 F15.922
F33.41 302.75 315.35 307.80 312.12 292.82 295.6 F06.2 F49 F86 F19.281 F11.251 F19.221 R40.2353 309. F15.188 312.4 F52.0 F19.19 F13.24 295.13 R42 308.9 F65.4
R45.6 309.9 307. F45.8 F11.122 R41.1 F30.3 F16.159 F32.89 F34.81 F14.981 F13.11 F10.929 F33.2 299.0 296.53 312.34 F16.10 F33.9 F65.51 R40.213 F17.200
F17.209 F12.20 F19.180 F10.180 296.6 296.04 F19.950 F31.12 F70 F55.1 F15.90 F32.81 F16.959 300.16 F45 F65.2 F14.250 F34 F17.290 F95.9 R44 F17.221 F19.929
300.21 R40.2252 295.7 R40.2324 F95 312.33 F15.10 295.04 314.00 F13.230 F14.159 308.2 F40 F16.120 F94.8 F60.1 295.84 F18.188 F40.243 F68.11 F18.920 F13.950
310. F12.150 F51.03 F13.129 F12.122 300.2 312.20 F16.988 295.93 F60.4 R46 296.31 F48.9 F07.9 R40.2321 298.9 F11.959 F31.63 F15.150 F16.221 F15.288
296.12 F16.983 F94 F11.120 F10.981 295.42 295.85 290.20 296.23 F48.8 R40.2420 F14.280 F30.2 R46.0 F19.150 F13.981 295.12 312.89 F15.222 F65.89 298.0
F19.922 R41.3 R40.2250 F15.181 F34.8 F57 F98.21 313.0 300.12 95.50 F25.9 F16.950 R40.2433 292.1 312.21 F14 307.46 F51 R40.2350 R49 F43.12 F91.0 306.
F17.229 F17.228 R40.2221 F43.0 F31.10 F14.11 R40.2352 F30.9 F40.2253 F14.90 R48.9 294.21 F12.11 R49.21 F14.229 293.81 296.32 F91.3 290
296.21 F11.259 F14.222 R40.20 F14.150 F14.921 296.51 307.42 F12.180 F15.121 R49.0 F18 292.9 300.19 F61 G44.209 F11.99 R45.7 309.89 F31.81 F19.280 295.72
F40.02 F22 F51.04 F30.13 F10.281 292.11 296.5 294 F12.19 F11.23 F13.231 F40.220 F14.950 F14.24 295.64 296.8 F40.230 300.10 F18.94 F16.951 F13.121 F19.159
F15.981 F12.10 313.82 296.34 F56 R47.81 F14.23 F13.159 F12.988 299.89 F19.288 295.74 300.9 R40.2320 F44 F31.64 296.41 F44.9
F19.982 F80.81 R40.2230 F13.281 R40.2354 F52.21 F34.0 R40.2144 F10.99 302.82 F29 295.10 F11.229 292.85 R40.2423 F12.921 300.5 F48.2 313.21 F10.182 310.81
F13.250 295.65 313.83 312.32 R48.3 R40.233 F31.75 F63.1 307.43 F50.01 F16.183 295.14 R40.2364 F88 302.83 F15.151 315.01 312.8 R41.844 F12.121 296.22 292.0
306.53 R41.82 R41.841 F53 R40.2123 313.9 312.22 296.60 R40.2312 300.0 302.52 R40.2120 307.20 R40.2143 F43.8 F18.99 F12.129 F51.9 F37 295.5 F01.50 296.7 R40.243 301.22 F12.280
F89 F60.9 F14.929 F60.6 F94.2 R40.2143 F43.8 F18.99 F12.129 F51.9 F37 295.5 F01.50 296.7 R40.243 301.22 F12.280

## Reproductive

764.96 O41.91x O70.0 628.0 646.03 O31.30x 634.91 608.22 679.10 662.11 O35.0xx0 P24.21 O13.9 656.01 'O34.00 P57.0 649.7 662.00 669.82 656.23 663.60 618.0
670.2 O34.529 671.80 O03.32 646.51 771.8 E28.2 668.03 674.80 651.2 653.50 661.11 629.9 N82.5 676.32 661.21 646.93 763.9 O36.0110 659.90 673.34 649.03 660.33
669.0 674.02 656.13 O92.5 675.11 668 659.10 674.50 N89.4 P07.24 669.50 660.20 617.1 602.2 O00.1 O99.350 N88.0 A48.51 660.31 O22.91 611.72 O01.9 651.50
661.2 648.51 642.53 634.50 O98.03 604.0 679.14 O35.4xx0 665.24 651.71 P35.1 655.70 O99.215 676 P91.63 773.4 767.0 616.51 671.10 O91.219 653.3 664.44 O72.0
669.94 670.8 761.7 673.11 P51.8 660.03 658.40 O00.0 651.7 765.27 671.90 O99.815 663.2 663.01 771.3 302.73 647.8 665.61 N70.93 611.0 O34.41 'O34.01 N64.2
631.0* 652.73 642.23 762 666.3 664.80 763 764.92 P05.06 N85.4 652 618.00 O36.61x0 O03.4 671.04 674.32 642.54 673.22 O75.89 N64.89 625.5 O48.0 'O33.5xx0
765.10 652.40 O42.00 764.12 O91.22 679.01 608.24 653.90 647.44 O41.91x0 647.84 643.81 630.0 762.4 659.30 656.4 651.20 669.6 O72.2 654.62 650.0 N64.3 649.60
O71.9 642.13 648.81 659.91 771.1 P01.0 654.93 603.9 654.14 670.80 674.52 768.1 O26.21 O33.5xx0 665.54 661.41 E23.0 602.1 'O33.9 652.50 O92.111 P03.0 O87.4
N88.8 642.92 N60.19 614.4 653.43 648.94 O99.330 O22.31 644.2 659.50 O71.3 646.3 760.70 651.21 'O33.6x 669.8 651.11 652.13 O91.12 762.6 614 660.43 O22.40
O32.6xx 760.61 653.4 E29.1 765.08 765.05 679 P00.1 656.93 647.01 P24.10 760.76 669.24 648.0 N48.6 N43.1 643.8 659.53 O69.2xx 647.10 651.23 N94.819 629.1
764.10 649.41 641.13 655.31 765.18 774.31 642.63 671.51 659.91 672.00 607 653.20 661.00 644.0 O41.1090 O35.4xx O41.90x 625.6 607.85 639.2 O99.419
O69.1xx 671.23 671.1 N81.2 664.8 O11.9 Q83.8 O99.341 665.60 604 656.31 661.23 655.33 669.00 656.7 653.70 614.7 646.8 661 646.14 620.1 645.03 O92.3 648.02
765.16 625.71 663.6 655.9 664.00 N71.0 651.31 768 O99.845 620.8 626.1 P35.2 656.41 652.90 646.54 P08.1 648.93 O22.00 662.1 655.1 N82.8 649.62 654.90 607.82
N80.9 O98.619 625.0 674.84 O31.31x0 R10.2 618 O41.8X90 647.20 653.80 679.02 648.52 O10.111 768.0 606.9 668.11 656.30 629.89 669.10 N89.6 670.10 642.93
N80.3 642.1 628 O86.0 660.93 630 646.01 N80.5 655.40 647.63 652.23 N92.4 P07.21 N94.818 656.71 765.02 O89.1 P03.6 O63.1 764.2 P56.0 N85.6 646.10 O23.91
654.03 670.32 654.33 P24.81 610.1 768.9 N81.12 648.73 652.10 O87.3 P50.3 656.54 662.30 761.3 O36.60x0 P03.89 669.1 P05.05 N94.810 658.01 O14.02 O60.12x0
N48.89 N95.2 O69.9xx 643.23 663.83 O88.23 673.0 771.0 626.5 665.6 659.70 662.23 667.14 647.50 648.14 P59.8 647.60 673.30 661.0 651.00 646.82 N83.4 674.90
669.51 620.4 765.07 670.20 646.1 639.5 P04.8 659.13 653.33 664.01 660.90 666.12 O41.8X1 N49.9 N71.9 653.71 N51 O26.879 O22.10 P22.0 624.5 665.0 P52.21
654.01 O99.03 647.61 665.9 634.31 669.91 O35.8xx0 O08.0 P24.30 659.33 N81.6 646.0 P24.31 675.04 658.21 674.10 N92.1 676.33 O88.019 669.04 676.1 668.23
668.80 665.81 762.8 654.54 763.4 O90.4 653.91 O25.2 670.84 647.64 766.21 674.3 659.60 P05.17 648.83 676.60 665.30 651.03 647.32 651.0 659.01 668.83 N46.9
767.5 610.0 668.2 675.23 661.40 642.33 768.3 644.03 N76.0 622.9 O92.6 O36.90x O88.311 654.00 656.00 601.8 668.04 P91.62 O99.340 761.4 765.03 659.3 652.6
648.43 674.9 669.43 643.93 O36.111 647.31 F52.32 647.3 O92.011 642.03 664.40 P11.3 608.8 668.22 641.80 O90.81 659.23 641.30 667 663.23 658.80 608.86 626.9
O66.1 665.64 O88.119 N64.1 P07.18 676.81 N92.2 N48.9 O03.7 O69.2xx0 O99.13 O08.6 657.0 673.00 664.81 N73.3 643.01 633.10 641.9 634.62 647.9 622 O99.280
N73.9 649.30 674.20 626.7 671.40 669.83 N41.9 651.6 647.04 764.04 628.39 764.14 O41.1010 O60.12x 665.40 674.6 46.83 O03.80 765.15 654.9
664.54 633.80 645.2 P39.0 616.1 O03.31 667.1 765 N64.4 649.11 N93.0 676.50 651.5 648.40 649.53 634.82 O99.345 O32.1xx0 665.83 O91.011 O30.029 653.61
O40.1xx0 O36.819 676.01 O35.9xx0 657.01 660.53 O03.83 647.43 663.43 P00.9 N43.3 774.39 773.2 664.11 653.10 608.3 623.6 O22.8X1 631 634.42 658.4 N83.9
654.84 774.30 639.8 768.5 651.73 763.84 766.1 607.1 648.80 O12.01 N46.029 669.40 O99.111 N41.8 773.1 763.89 602.8 633.8 618.7 674.4 621.6 O35.9xx 643.00
O26.849 634 664.04 764.23 611.4 O61.0 646.84 764.09 654.43 O35.3xx0 O31.11x 620.3 P01.7 656.83 E28.8 634.32 665.14 655.91 658.03 658.91 664.84 665.44
O91.02 611.5 647.62 607.2 665.3 N60.09 654.4 764.98 P03.819 658.10 666.10 O71.2 O66.5 N44.00 664.3 E28.0 O09.40 'O33.8 668.20 646.11 O71.02 642.60 656.51
678.10 P07.01 N85.3 659.4 O86.81 648.70 N94.9 O36.011 O82 614.1 652.01 O43.101 767.6 774.1 666.22 629 662.01 P01.2 O42.10 654.51 O36.1190 642.24 771.5
763.2 764.03 617.5 652.93 633.00 760.63 O66.9 O34.511 665.70 669.03 651.30 O33.3xx 642.4 660.10 649.12 675.80 N44.02 608.83 O36.829 602 O03.34 646.64
645.01 N90.9 616.4 653.6 N81.9 O36.91x 670.0 676.52 671.54 660.01 770.15 640.91 653.5 665.1 666.2 648.54 760.77 676.93 N61.1 P01.8 675.8 644.20 O34.21
656.33 765.11 618.83 608.9 627.3 P35.0 654.63 678.1 647.53 647.83 671.11 652.2 765.14 653.9 760.2 O43.019 611.71 P05.14 666.1 617.8 671.02 O91.23 674.0 O86.4
665.82 665.11 N75.0 651.83 O64.0xx 676.61 655.01 764.07 O26.619 O32.4xx0 O40.1xx 768.4 O26.41 671.82 671.9 648.1 601.4 641.21 P08.0 656.63 654.74 O76
652.41 675.20 653.00 O03.0 676.42 658.90 O98.019 642.43 659.20 661.90 660.73 658.9 641.01 641.00 676.51 N43.2 P24.80 675.6 60.0 O41.1010 625.9 618.04 664.14
N88.3 647.13 O65.5 666.14 648.72 676.11 O90.9 607.83 648.9 763.1 661.93 642.34 P02.4 653.03 760.4 655.7 674.2 O22.11 R78.81 764.00 671.2 661.10 P12.2
O30.201 646.9 667.12 765.09 633.0 770.86 620.7 O24.911 679.11 646.63 640.03 653.53 653.81 762.1 O99.019 P00.7 654.10 766.0 773.0 668.90 761.6 676.53 649.42
646.20 656.2 658 P05.18 651.91 O75.3 664.51 768.7 P07.00 P58.8 614.3 666.04 652.11 634.3 648.60 659.43 E29.0 653.11 N92.3 O03.39 667.02 649.00 646.42
E28.310 P04.0 O34.40 640.8 642.31 O22.51 641.81 O98.511 P02.9 675.92 764.24 626 P05.10 647.6 641.0 643.90 760.64 603.8 P05.13 639.4 620.6 667.00 O98.911
N97.2 P05.16 674.04 621.0 649.10 664.9 663.33 623.9 761.8 P10.3 F52.31 N94.1 346.4 624.0 P24.00 607.89 O25.10 O31.30x0 O65.9 O64.9xx 646.12 673.02 O36.8910
O90.5 767.19 670 764.17 652.60 764.99 642.90 634.1 N90.4 N90.89 648.8 618.1 633.01 646.13 642.5 664.60 641.33 O66.40 646.81 764 669.70 O25.3 O69.0xx0 665.01
646.60 O98.119 659.80 674.92 760.9 772.6 O90.3 664.61 621.1 663.61 608.20 774.4 P13.0 N95.8 N80.0 634.12

| Category | Codes |
|---|---|
| Respiratory | J27 J63.0 J12.9 J37.0 P26.0 J39 R05 J69 J15.1 J67.0 516.0 P24.21 J95.859 770.3 502 J09.X9 J34.3 J52 770.81 J03.00 786.9 J95.01 277.09 J20.6 J21.8 J33.8 P28.0 J95.850 J84.848 770.0 J67.7 P28.2 R06.7 J34.89 P27.8 478.30 J01.21 J12.3 518.3 G47.35 495.7 J11.82 J95.71 J64 J82 P28.11 J01.20 J15.8 J45.21 495.4 J39.9 J95.89 P23 P26.9 J95.811 J70.1 J73 J95.862 J66.1 J79 J93.83 J20.9 512 J41.8 J84.112 J03 518 518.1 J94.9 471.8 J18.1 J34.81 786.7 478.33 J34.2 J93 J32.9 R06 J68.8 J45.30 519.9 327.29 J20.0 J63.5 J84.03 P91.63 770.15 R06.81 495 J19.0 J01.10 770.7 J95.863 J09.X1 J02.0 J34.9 J84.02 J15.212 J81.0 J94.0 478.3 J23 J10.81 J84.116 495.1 507.0 J05.10 P26.1 J11.00 J06.0 J96.91 J33 R06.6 J45.991 516.32 770.12 491 J32.0 768.4 J87 768.70 J01.00 R06.01 J94.1 512.1 786.51 786.6 J38.02 J95.830 J38.01 478.79 770.84 J68.3 J10.08 J84.115 J25 J01.41 M34.81 P27.0 J07 512.8 R06.00 J11.08 J46 R09.3 J49 J70.5 J95.4 J05.0 J95.02 J31.0 J95.821 P24.80 J91.8 768.71 J30.5 J38 J70.2 J63.4 J86 478.4 J10.00 J84.9 J96.20 J21 J30.2 J31 J01.01 J29 R07.9 J75 J15.7 327.22 J61 J93.11 J02.9 478.75 P28.5 491.2 J45.22 J96.11 J67.8 J85.0 J35.03 J96.21 J11.1 J35.01 505 491.1 P28.10 J15.3 J38.00 J67 492.8 R06.2 770.88 770.86 J15.20 J38.5 R07 J30.1 J95.62 J20.1 J22 512.0 G47.30 768.1 P23.0 J35.3 J95.88 P91.60 J95.860 J81 J43.8 J00 J68.0 J83 478.34 J18.9 J94.2 J70.8 J35.9 J33.0 R07.1 770 501.0 J43.2 478.74 J45.20 470 J45.40 786.06 327.23 P28.89 R06.9 J39.0 J98.9 J47.0 J98.3 P22.1 J39.8 R06.83 J65 P24.10 R07.0 J37 R07.89 517 503 786.52 J10 491.20 J95.72 J28 J09.X3 J70 P23.3 J01.80 J97 J12.89 J18.8 516 J67.3 495.9 J40 P24.01 E84.9 J10.01 500.0 J11.2 519 501 J93.0 J98.01 P22 495.8 J08 G47.31 518.8 327.20 J20.5 786.01 786.50 P24.00 J41 770.2 J84.843 J93.81 515.0 J95.03 J30.0 786.00 J92 491.21 515 J92.9 768 R06.89 J56 J01.11 770.5 P19.1 J70.3 516.3 J32 J41.1 J85.3 J02.8 J21.9 J10.89 J68.2 518.2 R06.3 P25.1 R06.5 J91 P25.2 G47.36 J45.901 J16.8 J84.83 P28.4 R07.82 J18.2 J34 J74 J30.89 478.7 768.0 J11.83 J95.84 J59 478.5 J15.5 J84.17 J35.02 J18.0 768.72 786 519.02 J95.04 J63.6 P28.9 P19.9 J95 J85.1 491.8 J66.2 517.8 J95.822 J01.40 J48 519.09 518.4 G47.39 J51 786.59 277.03 J84.111 J86.0 P23.9 J34.0 J18 J17 J43 786.2 J96 J96.22 J71 J01.31 J63.3 495.6 494.0 P24.81 J95.2 J42 J76 J95.09 J84.09 768.9 J35.8 J01.91 G47.34 J06.9 J62.0 J96.10 J20.4 J50 491.9 R06.82 J32.2 J03.80 516.36 770.11 P23.4 J01.30 J84 J84.89 J47.1 J98.4 J01.90 R06.03 491.22 786.8 R07.81 492.0 J14 J66 P28.81 J15.211 J09.X2 J93.9 J67.5 J70.4 J45.32 84.09 770.6 J89 J01.81 P28.19 518.5 J32.8 J33.1 J10.82 786.05 478.6 J32.4 J81.1 J98.8 327.26 J96.02 516.8 518.0 519.4 770.9 J38.3 R06.09 J03.91 G47.37 327.2 J34.1 J12.0 J21.1 J01 P22.0 J12.2 J39.3 P25.3 786.02 J95.851 J31.1 P25.0 J15 J43.0 P23.1 495.2 478.31 P24.30 770.18 P24.31 J57 J63 J95.00 J69.0 J38.4 J09 277.0 J04.30 J95.5 J45.902 504 P24.11 J45 J03.01 P23.6 R06.1 519.0 786.1 J37.1 J15.29 J98.6 J44.9 J84.10 P27.9 J44.0 J96.00 517.3 E84.0 J45.998 516.9 J86.9 471 J63.2 491.0 J68.1 786.07 J98.19 P25 J66.8 768.3 518.84 J95.3 J30 J72 J44.1 519.00 J16 J84.2 770.85 J66.0 J20.2 519.11 277.00 J47.9 P27.1 770.4 J06 770.87 P91.62 P22.9 786.03 J84.117 J84.113 516.30 770.16 P23.2 495.0 J67.6 J93.82 J98.51 516.37 770.89 768.6 510 J04.2 J44 J12 327.24 'J43.9' J77 J84.82 J20.8 510.9 770.13 P28 J31.2 478.32 478.70 496 J38.6 J41.0 J24 J68 J11.89 J45.909 J03.81 J98.09 518.89 P91.61 770.10 J67.2 J15.0 J39.2 494.1 516.2 G47.33 277.02 J38.2 J60 E84.11 R09.89 J47 277.01 J95.831 J35 P23.5 J94.8 P24 J04.0 J69.8 517.2 R07.2 495.5 786.09 J20 J12.81 J45.31 J58 J45.51 J43.1 E84.19 P84 494 J55 J33.9 516.31 J30.9 P22.8 J98.59 J04.10 J38.1 J53 P25.8 P23.8 J62 516.33 J84.842 J54 J92.0 J32.1 769 J04.11 786.04 J68.9 J70.0 327.21 J90 P24.20 G47.32 770.83 P24.9 J12.1 J84.81 J67.1 J10.83 J32.3 J10.1 492 514 768.2 327.27 J95.61 J38.7 J16.0 J20.7 478.0 471.1 R06.02 J15.6 J95.1 J02 J95.812 J15.4 J68.4 J96.01 514.0 J03.90 J35.1 516.35 J45.990 786.3 J35.2 J62.8 J91.0 500 J93.12 770.82 J84.01 516.1 J45.41 J88 327.25 J84.114 P27 J20.3 471.9 J96.90 J11.81 518.82 495.3 J39.1 J98.11 J10.2 J80 519.8 J45.52 J96.12 518.83 J05.11 R06.4 770.17 J69.1 J67.4 J84.841 J99 J98.2 P26.8 J11 J63.1 J78 768.5 J85.2 J45.50 P26 P28.3 J15.9 J70.9 786.4 J36 J05 J95.861 J04.31 J67.9 J85 J04 478.71 471.0 518.81 E84.8 J98 J19 J21.0 J94 J43.9 770.14 516.34 J30.81 510.0 J45.42 J13 J26 768.73 R22.2 J96.92 |
| Rheumatism | 727.66 727.04 726.12 726.31 729.1 I20.8 726.72 728.89 729.89 728.83 727.49 729.4 I25.3 I21.11 727.06 726.10 727.43 729.73 I25.42 726.60 727.81 726.90 726.73 727.50 729.0 729.39 728.6 727.51 729.6 I25.10 726.62 727.67 729.71 727.40 726.69 727.68 726.5 I25.811 728.2 I25.2 727.09 727.59 726.30 727.05 I21.29 726.2 727.83 I21.4 726.39 728.85 727.60 729.92 728.9 726.4 726.91 727.03 729.82 I25.83 727.89 728.5 726.64 726.32 726.61 728.4 727.69 I25.9 728.88 728.10 727.65 727.00 728.79 I24.1 728.87 I25.41 727.82 726.70 729.81 729.91 I24.0 727.42 727.63 725 727.9 I21.19 I20.0 729.99 728.82 I20.1 727.3 727.64 726.63 726.65 727.2 727.01 726.79 729.31 728.81 726.33 729.30 728.11 I24.8 726.11 I25.89 728.12 729.72 728.0 727.41 729.90 727.02 726.0 I25.810 729.5 728.13 729.2 728.71 I21.3 728.3 727.1 I25.82 729.79 726.19 728.84 726.8 728.86 727.61 728.19 727.62 I21.09 I25.812 726.71 |
| Sleep-Disorders | 327.43 327.13 G47.37 G47.54 327.33 327.09 G47.01 I20.8 327.40 G47.35 G47.8 327.02 327.42 I25.3 I21.11 G47.51 327.15 G47.33 G47.419 327.32 327.11 G47.10 I25.42 327.20 327.27 327.51 G47.13 G47.27 327.41 G47.19 I25.10 G47.9 G47.22 G47.39 I25.811 327.00 327.36 327.30 G47.61 I25.2 327.34 G47.31 G47.69 I21.29 327.49 G47.14 I21.4 G47.50 327.53 327.31 I25.83 G47.29 G47.30 327.25 G47.36 327.35 G47.09 I25.9 G47.32 G47.52 327.24 I24.1 I25.41 G47.63 G47.23 I24.0 327.01 G47.11 G47.421 G47.20 327.10 G47.34 G47.00 327.29 I21.19 I20.0 I20.1 327.21 327.44 G47.24 327.52 327.19 I24.8 327.8 I25.89 327.26 327.37 G47.12 G47.429 G47.53 327.22 I25.810 G47.26 G47.411 G47.62 G47.21 327.39 I21.3 327.12 I25.82 327.23 327.59 G47.25 327.14 I21.09 I25.812 G47.59 |
| Symptoms-Abs-Pelvis | R10.2 R10.11 R10.83 789.7 789.06 I20.8 789.02 789.1 R10.10 I25.3 I21.11 R10.812 789.42 789.63 789.04 789.2 R10.815 I25.42 R10.819 789.44 789.61 R10.811 789.34 789.47 R10.13 789.60 789.09 I25.10 789.67 I25.811 789.69 789.65 789.33 I25.2 789.01 I21.29 789.35 I21.4 R10.31 789.51 R10.12 I25.83 789.49 789.05 I25.9 R10.817 789.64 I24.1 I25.41 789.36 I24.0 R10.33 789.9 I21.19 I20.0 I20.1 R10.813 789.45 789.32 789.62 R10.9 789.43 789.40 I24.8 I25.89 789.46 789.41 R10.32 I25.810 789.00 I21.3 R10.84 789.31 I25.82 789.37 789.30 789.03 789.66 789.39 R10.814 R10.816 789.07 I21.09 I25.812 789.59 |
| Symptoms-Digestive | R14.1 R13.14 I20.8 787.5 787.91 R19.4 R19.05 R19.01 R19.2 I25.3 I21.11 R11.11 R11.2 R19.07 R19.06 R19.09 R15.0 I25.42 787.3 R19.32 R13.13 R19.03 787.4 R19.30 787.21 R11.13 R15.9 I25.10 787.04 I25.811 I25.2 R11.0 787.20 I21.29 R19.35 R11.10 I21.4 R19.33 R19.04 787.29 787.03 I25.83 R16.0 R19.31 R12 R13.19 R17 R19.02 R15.1 R18.0 787.6 I25.9 R16.1 I24.1 I25.41 787.02 R15.2 I24.0 R13.12 R11.14 787.99 R13.11 787.23 I21.19 I20.0 I20.1 R19.36 R19.11 R19.8 I24.8 R19.37 I25.89 R13.10 R18.8 787.7 787.24 R19.34 787.01 I25.810 787.22 I21.3 R19.7 I25.82 R19.5 R19.00 I21.09 I25.812 787.1 |
| Symptoms-General | R53.82 R65.10 R56.01 R50.81 R68.3 R59.9 780.39 I20.8 R53.1 R62.0 R55 R65.21 R62.51 780.64 I25.3 I21.11 780.1 780.63 780.65 R56.1 R68.0 R62.50 R64 780.59 R63.3 I25.42 780.96 R68.83 780.8 R56.00 R68.84 R62.52 R53.81 I25.10 R63.0 780.72 780.51 R68.82 780.91 R63.1 I25.811 I25.2 R68.12 I21.29 R63.2 R57.0 I21.4 R63.6 R63.4 R63.8 780.92 R57.9 I25.83 780.62 780.52 R65.20 R65.11 780.55 R68.81 R50.82 780.32 780.94 I25.9 780.95 780.71 I24.1 R50.84 I25.41 780.56 R68.11 780.09 I24.0 R50.9 R63.5 I21.19 I20.0 I20.1 780.99 780.97 780.61 R53.2 780.4 I24.8 R58 R69 780.2 R51 I25.89 R56.9 780.93 R52 780.03 R60.9 780.58 780.01 780.53 I25.810 780.31 780.57 R62.7 I21.3 I25.82 R61 780.79 780.50 R57.8 R68.13 780.02 780.54 780.60 I21.09 I25.812 R50.83 R68.89 |

| | |
|---|---|
| Symptoms-Respiratory | 786.7 786.50 R09.2 786.3 R06.4 786.51 R06.7 R09.82 I24.8 I25.10 786.04 786.8 R06.9 786.2 I25.89 I20.8 R07.0 786.02 I25.9 786.06 786.05 786.1 I25.811 R06.3 R06.2 R09.02 I24.1 I25.2 786.4 I25.41 I25.810 I25.3 I21.11 R06.81 R06.02 R06.1 786.07 R06.01 R06.6 I21.29 I21.3 786.00 786.6 786.9 I24.0 R05 I25.82 I21.4 R06.89 I25.42 R06.82 R09.3 R06.00 786.09 786.59 R09.89 R09.01 I21.19 I20.0 R07.2 R07.1 I20.1 I21.09 I25.812 786.03 R07.89 786.52 786.01 I25.83 R07.9 |
| Symptoms-Skin | R22.1 R23.4 R23.8 I24.8 I25.10 I25.89 I20.8 R23.2 782.2 I25.9 R22.2 I25.811 782.4 I24.1 I25.2 782.0 I25.41 I25.810 I25.3 I21.11 782.1 782.62 R21 R23.0 782.7 782.61 I21.29 I21.3 I24.0 R23.3 782.3 782.8 I25.82 I21.4 I25.42 782.5 I21.19 I20.0 I20.1 I21.09 I25.812 R23.1 R22.9 R20.3 782.9 I25.83 |
| Symptoms-Urinary | R36.1 R36.9 788.35 788.61 788.36 I20.8 788.0 788.63 I25.3 I21.11 788.69 I25.42 R37 788.32 788.20 788.39 R32 R39.12 I25.10 R39.14 R31.9 788.65 788.38 R39.81 R35.8 R39.16 I25.811 R33.9 I25.2 I21.29 788.30 I21.4 788.37 R39.0 788.43 R39.13 R39.89 788.31 I25.83 788.99 788.7 R39.11 I25.9 I24.1 I25.41 R39.19 I24.0 788.91 I21.19 I20.0 788.8 R35.0 I20.1 788.21 788.33 788.42 I24.8 R35.1 R30.0 I25.89 R34 R33.8 788.62 I25.810 788.41 R39.15 I21.3 788.1 788.34 788.64 I25.82 788.5 788.29 I21.09 I25.812 R31.0 R31.1 |
| Thyroid | E03.9 241.9 E04.1 E04.2 E05.90 I20.8 E05.21 I25.3 I21.11 245.2 243 241.0 E05.40 I25.42 E05.91 246.9 E07.89 E04.0 242.10 242.30 245.0 E07.0 I25.10 242.81 242.11 244.3 242.41 242.40 246.2 I25.811 240.0 E06.4 I25.2 I21.29 I21.4 242.21 241.1 E05.11 I25.83 E06.3 245.8 242.80 E05.20 E06.1 I25.9 E05.10 I24.1 I25.41 E07.81 E05.01 E01.8 245.4 I24.0 246.3 242.91 244.8 246.1 242.20 I21.19 I20.0 E01.2 242.01 I20.1 244.0 242.00 I24.8 245.1 E05.00 I25.89 244.9 246.0 E07.1 E06.9 I25.810 E00.9 244.1 E05.41 I21.3 244.2 242.90 240.9 245.9 E07.9 I25.82 E05.31 E05.30 245.3 E04.9 242.31 E06.0 246.8 E03.2 I21.09 I25.812 E06.5 |

---

**Algorithm 1:** `GenESeSS`

---

**Data:** A sequence $x$ over alphabet $\Sigma$, $0 < \varepsilon < 1$

**Result:** State set $Q$, transition map $\delta$, and transition probability $\widetilde{\pi}$

    /* **Step One: Approximate $\varepsilon$-synchronizing sequence**                                             */

1   Let $L = \left\lceil \log_{|\Sigma|} 1/\varepsilon \right\rceil$;

2   Calculate the **derivative heap** $\mathcal{D}_\varepsilon^x$ equaling $\left\{ \hat{\phi}_y^x \; : \; y \text{ is a sub-sequence of } x \text{ with } |y| \le L \right\}$;

3   Let $\mathcal{C}$ be the convex hull of $D_\varepsilon^x$;

4   Select $x_0$ with $\hat{\phi}_{x_0}^x$ being a vertex of $\mathcal{C}$ and has the highest frequency in $x$;

    /* **Step Two: Identify transition structure**                                                      */

5   Initialize $Q = \{q_0\}$;

6   Associate to $q_0$ the **sequence identifier** $x_{q_0}^{\mathsf{id}} = x_0$ and the probability vector $d_{q_0} = \hat{\phi}_{x_0}^x$;

7   Let $\widetilde{Q}$ be the set of states that are just added and initialize it to be $Q$;

8   **while** $\widetilde{Q} \ne \emptyset$ **do**

9      Let $Q_{\mathsf{new}} = \emptyset$ be the set of new states;

10      **for** $(q, \sigma) \in \widetilde{Q} \times \Sigma$ **do**

11          Let $x = x_q^{\mathsf{id}}$ and $d = \hat{\phi}_{x\sigma}^x$;

12          **if** $\|d - d_{q'}\|_\infty < \varepsilon$ *for some* $q' \in Q$ **then**

13              Let $\delta(q, \sigma) = q'$;

14          **else**

15              Let $Q_{\mathsf{new}} = Q_{\mathsf{new}} \cup \{q_{\mathsf{new}}\}$ and $Q = Q \cup \{q_{\mathsf{new}}\}$;

16              Associate to $q_{\mathsf{new}}$ the sequence identifier $x_{q_{\mathsf{new}}}^{\mathsf{id}} = x\sigma$ and the probability vector $d_{q_{\mathsf{new}}} = d$;

17              Let $\delta(q, \sigma) = q_{\mathsf{new}}$;

18      Let $\widetilde{Q} = Q_{\mathsf{new}}$;

19   Take a strongly connected subgraph of the labeled directed graph defined by $Q$ and $\delta$, and denote the vertex set of the subgraph again by $Q$;

    /* **Step Three: Identify transition probability**                                       */

20   Initialize counter $N[q, \sigma]$ for each pair $(q, \sigma) \in Q \times \Sigma$;

21   Choose a random starting state $q \in Q$;

22   **for** $\sigma \in x$ **do**

23      Let $N[q, \sigma] = N[q, \sigma] + 1$;

24      Let $q = \delta(q, \sigma)$;

25   Let $\widetilde{\pi}(q) = \left\| (N[q, \sigma])_{\sigma \in \Sigma} \right\|$;

26   **return** $Q, \delta, \widetilde{\pi}$;

---

**Algorithm 2:** Log-likelihood

---

**Data:** A PFSA $G = (\Sigma, Q, \delta, \widetilde{\pi})$ and a sequence $x$ over alphabet $\Sigma$

**Result:** Log-likelihood $L(x, G)$ of $G$ generating $x$

1   Calculate the state transition matrix $\Pi$ and observation $\widetilde{\Pi}$;

2   Calculate the stationary distribution over states $\wp_G$ of $G$ from $\Pi$;

3   Calculate the stationary distribution of alphabet $\phi_\lambda^T = \wp_G^T \widetilde{\Pi}$;

4   Initialize $\mathbf{p}$ by $\wp_G$ and $\mathbf{q}$ by $\phi_\lambda$;

5   Let $L = 0$;

6   **for** $i$ *from* $1$ *to* $|x|$ **do**

7      Let $\sigma$ be the $i$-th entry of $x$;

8      Let $L = L - \log \mathbf{q}|_\sigma$;

9      Let $\mathbf{p}^T = \left\| \mathbf{p}^T \Gamma_\sigma \right\|$ where $\Gamma_\sigma$ is defined in Eq. 8;

10      Let $\mathbf{q}^T = \mathbf{p}^T \widetilde{\Pi}$;

11   **return** $L/|x|$;