

Reduced False Positives in Autism Screening Via Digital Bio-markers Inferred from Deep Co-morbidity Patterns

Dmytro Onishchenko^a, Yi Huang^a, James van Horne^a, Peter J. Smith^{d,g}, Michael M. Msall^{e,f}, and Ishanu Chattopadhyay^{a,b,c,1}

^aDepartment of Medicine, University of Chicago; ^bCommittee on Genetics, Genomics & Systems Biology, University of Chicago; ^cCommittee on Quantitative Methods in Social, Behavioral, and Health Sciences, University of Chicago; ^dDepartment of Pediatrics, Section of Developmental and Behavioral Pediatrics, University of Chicago; ^eDepartment of Pediatrics, Section Chief of Developmental and Behavioral Pediatrics, University of Chicago; ^fJoseph P. Kennedy Research Center on Intellectual and Neurodevelopmental Disabilities, University of Chicago; ^gExecutive Committee Chair, American Academy of Pediatrics' Section on Developmental and Behavioral Pediatrics

This manuscript was compiled on Friday 19th June, 2020

1 **Autism spectrum disorder (ASD) is a developmental disability associated with significant social, communication, and behavioral challenges. There is a need for tools that help identify children with ASD as early as possible (1, 2). Our current incomplete understanding of ASD pathogenesis, and the lack of reliable biomarkers hampers early detection, intervention, and developmental trajectories. In this study we develop and validate machine inferred digital biomarkers for autism using individual diagnostic codes already recorded during medical encounters. Our risk estimator identifies children at high risk with a corresponding area under the receiver operating characteristic curve (AUC) exceeding 80% from shortly after two years of age for either sex, and across two independent databases of patient records. Thus, we systematically leverage ASD co-morbidities - with no requirement of additional blood work, tests or procedures - to predict elevated risk during the earliest childhood years, when intervention is the most effective. Our methodology has superior performance to common questionnaires-based screenings such as the M-CHAT/F (3), and has the potential to reduce socio-economic, ethnic and demographic biases. Further, by conditioning on the individual M-CHAT/F scores, we can either halve the false positives or boost sensitivity by over 50%, while maintaining specificity above 95%. Translated into practice, our algorithmic approach could significantly reduce the median diagnostic age for ASD, and also reduce long post-screen wait-times (4) currently experienced by families for confirmatory diagnoses and access to evidence based interventions.**

autism | electronic health record | stochastic learning | co-morbid risk

1 **A**utism spectrum disorder is a developmental disability associated with significant social, and behavioral challenges. Even though ASD may be diagnosed as early as the age of two (5), children frequently remain undiagnosed until after the fourth birthday (6). At this time, there are no laboratory tests for ASD, so a careful review of behavioral history, and a direct observation of symptoms is necessary (7, 8) for a clinical diagnosis. Starting with a positive initial screen, a confirmed diagnosis of ASD is a multi-step process that often takes 3 months to 1 year, delaying entry into time-critical intervention programs. While lengthy evaluations (9), cost of care (10), lack of providers (11), and lack of comfort in diagnosing ASD by primary care providers (11) are all responsible to varying degrees (12), one obvious source of this delay is the number of false positives produced in the initial ASD-specific screening tools in use today. For example, the M-CHAT/F used commonly as a screening tool (8, 13), has an estimated sensitivity of 38.8%, specificity of 94.9% and Positive Predictive Value (PPV) of 14.6% (3). Thus, currently out of every 100 children with ASD, M-CHAT/F flags about 39, and out of every 100 children it flags, about 85 are false positives, exacerbating wait

times and queues (12). Automated screening that might be administered with no specialized training, requires no behavioral observations, and is functionally independent of the tools employed in current practice, has the potential for immediate transformative impact on patient care.

While the neurobiological basis of autism remains poorly understood, a detailed assessment conducted by the US Centers for Disease Control and Prevention (CDC) demonstrated that children with ASD experience higher than expected rates of many diseases (5). These include conditions related to dysregulation of immune pathways such as eczema, allergies, asthma, as well as ear and respiratory infections, gastrointestinal problems, developmental issues, severe headaches, migraines, and seizures (14, 15). In the present study, we exploit these comorbidities to estimate the risk of childhood neuropsychiatric disorders on the autism spectrum. Using sequences of diagnostic codes from past doctor's visits, our risk estimator reliably predicts an eventual clinical diagnosis | or the lack thereof | for individual patients. Thus, the key clinical contribution of this study is the formalization of subtle co-morbidity patterns as a reliable screening tool, and potentially improve wait-times for diagnostic evaluations by significantly reducing the number of false positives encountered in initial screens in current practice.

A screening tool that tracks the risk of an eventual ASD diagnosis, based on the information already being gathered during regular doctor's visits, and which may be implemented as a fully automated background process requiring no time commitment from providers has the potential to reduce avoid-

Significance Statement

In this study past medical history in the form of diagnostic codes generated during medical encounters is used to formulate a clinically useful risk estimator for autism screening. We can formulate such a predictor despite considerable heterogeneity in ASD presentation that has made reliable biomarkers hard to identify, and has implicated over 1000 genes as contributing to ASD risk. Furthermore, the ability to carry out this evaluation without any blood tests, or questionnaires, makes it possible to avoid uncertainties from language barriers and other socio-economic artifacts. In joint operation with existing screening tools we can boost the predictor performance significantly higher than the current state of art, potentially cutting down false positives by half without losing specificity.

DO implemented the algorithm, DO, YH, JH and IC worked on mathematical modeling, IC, PS and MM designed study parameters, IC wrote the paper

²To whom correspondence should be addressed. E-mail: ishanu@uchicago.edu

Table 1. Patient Counts In De-identified Data & The Fraction of Datasets Excluded By Our Exclusion Criteria*

Distinct Patients	Truven		UCM	
	Male	Female	Male	Female
ASD Diagnosis Count†	12,146	3,018	307	70
Control Count†	2,301,952	2,186,468	20,249	17,386
AUC at 125 weeks	82.3%	82.5%	83.1%	81.37%
AUC at 150 weeks	84.79%	85.26%	82.15%	83.39%

Excluded Fraction of the Data sets				
Positive Category	0.0002	0.0	0.0160	0.0
Control Category	0.0045	0.0045	0.0413	0.0476

Average Number of Diagnostic Codes In Excluded Patients (corresponding number in included patients)				
Positive Category	4.33 (35.93)	0.0 (36.07)	2.6 (9.75)	0.0 (10.18)
Control Category	1.57 (17.06)	1.48 (15.96)	2.32 (6.8)	2.07 (6.79)

† Cohort sizes are smaller than the total number of distinct patients due to the following exclusion criteria: 1) At least one code within our complete set of tracked diagnostic codes is present in the patient record, 2) Time-lag between first and last available record for a patient is at least 15 weeks.

* Dataset sizes are after the exclusion criteria are applied

Table 3. Engineered Features (Total Count: 165)

Feature Type‡	Description	No. of Features
[Disease Category] Δ	Likelihood Defect (See Methods section)	17
[Disease Category] 0	Likelihood of control model (See Methods section)	17
[Disease Category] proportion	Occurrences in the encoded sequence / length of the sequence	17
[Disease Category] streak	Maximum Length of adjacent occurrences of [Disease Category]	51
[Disease Category] prevalence	Maximum, mean and variance of Occurrences in the encoded sequence / Total Number of diagnostic codes in the mapped sequence	51
Feature Mean, Feature Variance, Feature Maximum for difference of control and case models	Mean, Variance, Maximum of the [Disease Category] Δ values	3
Feature Mean, Feature Variance, Feature Maximum for control models	Mean, Variance, Maximum of the [Disease Category] 0 values	3
Streak	Maximum, mean and variance of the length of adjacent occurrences of [Disease Category]	3
Intermission	Maximum, mean and variance of the length of adjacent empty weeks	3

‡ Disease categories are described in Table 2

able diagnostic delays at no additional burden of time, money and personnel resources. While still lacking the certainty of a diagnostic blood test, use of patterns emergent in the diagnostic history to estimate risk might help reduce the subjective component in questionnaire-based screening tools, resulting in 1) reduced effect of potential language and cultural barriers in diverse populations, and 2) possibly better identify children with milder symptoms (8). Furthermore, being functionally independent of the M-CHAT/F, we show that there is clear advantage to combining the outcomes of the two tools: we can take advantage of any population stratification induced by the M-CHAT/F scores to significantly boost combined screening performance (See Materials & Methods, and Supplementary text, section 9).

Results

We measure our performance using several standard metrics including the AUC, sensitivity, specificity and the PPV. For

the prediction of the eventual ASD status, we achieve an out-of-sample AUC of 82.3% and 82.5% for males and females respectively at 125 weeks for the Truven dataset. In the UCM dataset, our performance is comparable: 83.1% and 81.3% for males and females respectively (Fig. 2 and 3). Our AUC is shown to improve approximately linearly with patient age: Fig. 3A illustrates that the AUC reaches 90% in the Truven dataset at the age of four. Importantly, we train our pipeline on 50% of the Truven dataset, and use held back data from Truven, and the entirety of the UCM dataset for validation: *No new training is done in the UCM dataset*. Good performance on these independent datasets lends strong evidence for our claims. Furthermore, applicability in new datasets *without local re-training* makes it readily deployable in clinical settings.

What are the inferred patterns that elevate risk? Enumerating the top 15 predictive features (Fig. 2B), ranked according to their automatically inferred weights (the feature “importances”), we found that while infections and immunologic disorders are the most predictive, there is significant effect

Table 2. Disease Categories (A few ICD9 codes shown from the complete set of 9,835 unique ICD9 codes considered. See SI-Table 4 in Supplementary text for complete list)

Category [†]	Description	Examples of ICD9* Codes
ASD*	Diagnostic Target	299 299.0 299.00 299.01 299.9 299.8 299.91 299.90 299.80 299.81 299.1 299.10 299.11
Immunologic	Diseases related to dys-regulation of the Immune system	580.81 580.89 580.0 580.8 461 461.8 461.0 477.9 477.2 477 477.8
Infectious	Diseases Caused By Pathogens	487.8 488.12 488.0 488.01 487.0 487.1 488.09 464.4 466 466.11 466.1
Nutrition	Symptoms concerning nutrition, metabolism and development	783.0 783.21 783.3 783.40 783.42 783.7 783.9
Mental Disorders	Psychiatric phenotypes other than ASD	290 - 319 (except 299.x)
Health Services	Contact With Health Services and Classification Of Factors Influencing Health Status	V01.0 V01.1 V01.2 V01.3 V01.4 V09.70 V09.71 V88.02 V88.03 V89.01 V89.02 V89.03 V89.04 V89.05 V89.09
Digestive	Diseases Of The Digestive System	540.0 540.1 541.0 542 540 541 543.0 562.03 562.01 562.00 562.10
Otic	Diseases Of The Ear And Mastoid Process	381.51 381.50 381.81 381.89 381.61 381.62 381 381.7 385.82 383.32 380.30
Musculoskeletal	Congenital musculoskeletal anomalies	756.52 756.53 733.02 733.0 733.09 737.43 737.41 737.20 737.29 737.4 737.2
Developmental	Congenital anomalies (Non-overlapping with musculoskeletal)	755.55 743.45 743.11 743.10 743.00 743.03 743.44 743.22 743.20 743.21 758.4
Reproductive	Diseases Of The Genitourinary System	611.79 611.71 611.89 611.81 676.64 611 676.60 611.6 611.4 611.3 611.2
Integumentary	Diseases Of Skin And Subcutaneous Tissue	706.0 706.1 704.00 704.02 704.09 680.9 680.1 680.5 680.7 680.6 680
Ophthalmologic	Disorders Of The Eye And Adnexa	362.8 362.9 362.6 362.1 362.3 362.18 362.17 362.13 362.11 363.33 363.32
Hematologic	Diseases Of The Blood And Blood-Forming Organs	286.9 286.6 283.19 283.11 283.9 283.1 284.0 284.09 284 284.01
Metabolic	Metabolic Disorders (Non-overlapping with respiratory, digestive and immunological conditions)	273.4 270 270.3 712.11 712.13 712.12 712.14 712.18 712.30 712.37 712.36
Cardiovascular	Diseases Of Arteries, Arterioles, And Capillaries	442.89 441.6 442.82 442.83 441.03 441.02 441.00 442 414.11 447.70 447.71
Respiratory	Diseases Of The Respiratory System (non-overlapping with Infectious)	516.31 516.30 516.32 516.35 516.37 516.36 516.8 516.0 277.0 277.00 277.01
Endocrine	Disorders Of Thyroid and other Endocrine Glands	244 244.9 244.2 255.41 255.5 255.4 259.51 255 259.4 255.11 242.2

[†] Categories inferred to be important for risk modulation are proportionately highlighted.

* ICD10 codes when present were mapped back to closest ICD9 matches using published General Equivalence Mappings (17).

from all the 17 disease categories. Thus, the co-morbid indicators are distributed across the disease spectrum, and no single disorder is uniquely implicated (See also Fig. 3F). Importantly, predictability is relatively agnostic to the number of local cases across US counties (Fig. 2C-D) which is important in light of the current uneven distribution of diagnostic resources (12, 18) across states and regions.

Unlike individual predictions which only become relevant over 2 years, the average risk over the populations is clearly different from around the first birthday (Fig. 3B), with the risk for the positive cohort rapidly rising. Also, we see a saturation of the risk after \approx 3 years, which corresponds to the median diagnosis age in the database (See Fig. 1B). Thus, if a child is not diagnosed up to that age, then the risk falls, since the probability of a diagnosis in the population starts to go down after this age. While average discrimination is not useful for individual patients, these reveal important clues as to how the risk evolves over time. Additionally, while each new diagnostic code within the first year of life increases the risk burden by approximately 2% irrespective of sex (Fig. 3D), distinct categories modulate the risk differently, *e.g.*, for a single random

patient illustrated in Fig. 3F infections and immunological disorders dominate early, while diseases of the nervous system and sensory organs, as well as ill-defined symptoms, dominate the latter period.

Given these results, it is important to ask how much earlier can we trigger an intervention? On average, the first time the relative risk (risk divided by the decision threshold set to maximize F1-score, see Methods) crosses the 90% threshold precedes diagnosis by \approx 188 weeks in the Truven dataset, and \approx 129 weeks in the UCM dataset. This does not mean that we are leading a possible clinical diagnosis by over 2 years; a significant portion of this delay arises from families waiting in queue for diagnostic evaluations. Nevertheless, since delays are rarely greater than one year (12), we are still likely to produce valid red flags significantly earlier than the current practice.

Our approach produces a strictly superior PPV (exceeding M-CHAT/F PPV by at least 14% (14.1-33.6%) when sensitivity and specificity are held at comparable values (approx. 38% and 95%) around the age of 26 months (\approx 112 weeks). Fig. 4A and Table 4 show the out-of-sample PPV vs sensitivity curves

Table 4. PPV Achieved at 100, 112 and 150 Weeks For Each Dataset and Gender (M-CHAT/F: sensitivity=38.8%, specificity=95%, PPV=14.6% between 16 and 26 months (\approx 112 weeks))

weeks	specificity	sensitivity	PPV	gender	dataset
100	0.92	0.39	0.14	F	UCM
100	0.95	0.39	0.19	M	UCM
100	0.93	0.39	0.13	F	Truven
100	0.91	0.39	0.10	M	Truven
112	0.93	0.39	0.16	F	UCM
112	0.95	0.39	0.20	M	UCM
112	0.96	0.39	0.22	F	Truven
112	0.95	0.39	0.17	M	Truven
150	0.94	0.39	0.19	F	UCM
150	0.98	0.39	0.34	F	Truven
150	0.97	0.39	0.26	M	Truven
150	0.97	0.39	0.26	M	UCM

Table 5. Personalized Operation Conditioned on M-CHAT/F Scores at 26 months

M-CHAT/F Outcome				global performance (Truven)			global performance (UCM)			prevalence*
0-2 NEG	3-7 NEG	3-7 POS	\geq 8 POS	speci-ficity	sensi-tivity	PPV	speci-ficity	sensi-tivity	PPV	
0.2	0.54	0.83	0.98	0.95	0.585	0.209	0.95	0.505	0.186	0.022
0.21	0.53	0.83	0.98	0.95	0.586	0.208	0.95	0.506	0.184	0.022
0.42	0.87	0.98	0.99	0.98	0.433	0.331	0.98	0.347	0.284	0.022
0.48	0.87	0.97	0.99	0.98	0.432	0.331	0.98	0.355	0.289	0.022
0.38	0.54	0.94	0.98	0.95	0.736	0.203	0.95	0.628	0.178	0.017
0.3	0.55	0.94	0.98	0.95	0.737	0.203	0.95	0.633	0.179	0.017
0.58	0.96	0.98	0.99	0.98	0.492	0.302	0.98	0.373	0.247	0.017
0.59	0.96	0.98	0.99	0.98	0.491	0.303	0.98	0.372	0.248	0.017
0.46	0.92	0.97	0.99	0.977	0.534	0.291	0.977	0.448	0.256	0.017
0.48	0.92	0.97	0.99	0.978	0.533	0.292	0.978	0.448	0.257	0.017

*Prevalence reported by CDC is 1.7%, while the CHOP study reports a value of 2.23%. The results of our optimization depend on the prevalence estimate.

for the two databases, stratified by sex, computed at 100, 112 and 100 weeks. A single illustrative operating point is also shown on the ROC curve in Fig. 2C, where at 150 weeks, we have a sensitivity of 51.8% and a PPV of 15.8% and 18.8% for males and females respectively, both at a specificity of 95%.

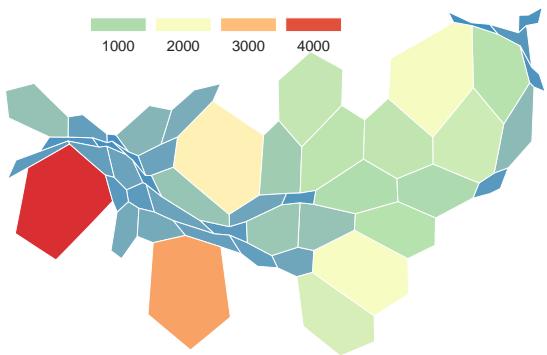
Beyond standalone performance, independence from standardized questionnaires implies that we stand to gain substantially from combined operation. With the recently reported population stratification induced by M-CHAT/F scores (3) (SI-Table 3), we can compute a conditional choice of sensitivity for our tool, in each sub-population (M-CHAT/F score brackets: 0 – 2, 3 – 7 (negative assessment), 3 – 7 (positive assessment), and $>$ 8), leading to a significant performance boost. With such conditional operation, we get a PPV close to or exceeding 30% at the high precision (HP) operating point across datasets ($>$ 33% for Truven, $>$ 28% for UCM), or a sensitivity close to or exceeding 50% for the high recall (HR) operating point ($>$ 58% for Truven, $>$ 50% for UCM), when we restrict specificities to above 95% (See Table 5, Fig. 4B(i), and SI-Fig. 4 in the supplementary text). Comparing with standalone M-CHAT/F performance (Fig. 4B(ii)), we show that for any prevalence between 1.7% and 2.23%, we can *double the PPV* without losing sensitivity at $>$ 98% specificity, or increase the sensitivity by \sim 50% without sacrificing PPV and keeping specificity \geq 94%.

Discussion

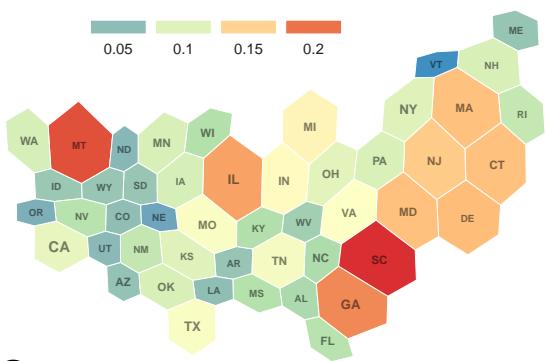
In this study, we operationalize a documented aspect of ASD symptomatology in that it has a wide range of co-morbidities (14, 15, 19) occurring at above-average rates (8). Association of ASD with epilepsy (20), gastrointestinal disorders(21–26), mental health disorders (27), insomnia, decreased motor skills (28), allergies including eczema (21–26), immunologic (19, 29–35) and metabolic(25, 36, 37) disorders are widely reported. These studies, along with support from large scale exome sequencing (38, 39), have linked the disorder to putative mechanisms of chronic neuroinflammation, implicating immune dysregulation and microglial activation (31, 34, 40–43) during important brain developmental periods of myelination and synaptogenesis. However, these advances have not yet led to clinically relevant diagnostic biomarkers. Majority of the co-morbid conditions are common in the control population, and rate differentials at the population level do not automatically yield individual risk (44).

Attempts at curating genetic biomarkers has also met with limited success. ASD genes exhibit extensive phenotypic variability, with identical variants associated with diverse individual outcomes not limited to ASD, including schizophrenia, intellectual disability, language impairment, epilepsy, neuropsychiatric disorders and, also typical development (45). Additionally, no single gene can be considered “causal” for

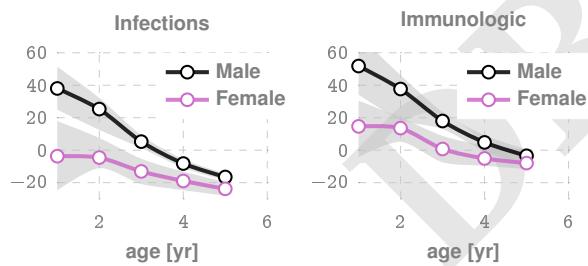
A. Autism Insurance Claims 2003-2013
(source: Truven Marketscan)



B. Autism Prevalence in US (Population Normalized)



C. Population-level Prevalence Differences between Positive vs Control Populations



D. ASD Clinical Diagnosis Age Across Genders

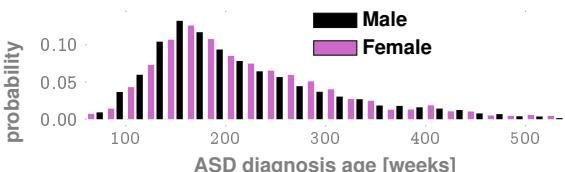


Fig. 1. ASD Occurrence Patterns. Panel A illustrates the differential representation of different disease categories in the positive and control cohorts, and panel B shows the distribution of the age of diagnosis for males and females in the Truven dataset. Panel C illustrates the spatial distribution of ASD insurance claims, and panel D shows the same data after population normalization, illustrating the relatively small demographic skew to ASD prevalence within the general population with access to medical insurance, which is consistent with the suggestion that prevalence variation might be linked to regional and socioeconomic disparities in access to services (16).

gorize behavior. This is susceptible to potential interpretative biases arising from language barriers, as well as social and cultural differences, often leading to systematic under-diagnosis in diverse populations (8). In this study we use time-stamped sequence of past disorders to elicit crucial information on the developing risk of an eventual diagnosis, and formulate a screening protocol that is free from such biases, and yet significantly outperforms the tools in current practice.

Going beyond screening performance, this approach provides a new tool to uncover clues to ASD pathobiology. Perhaps this vulnerability to diverse immunological, endocrinological and neurological impairments reflects how allostatic loads of medical stress get under the skin and disrupt key regulators of CNS organization and synaptogenesis. Charting individual disorders in the co-morbidity burden further reveals novel associations in normalized prevalence | the odds of experiencing a specific disorder, particularly in the early years (age < 3 years), normalized over all unique disorders experienced in the specified time-frame. We focus on the true positives in the positive cohort and the true negatives in the control cohort to investigate patterns that correctly disambiguate ASD status. On these lines Fig. 5 and SI-Fig. 1 in the supplementary text outline two key observations: 1) *negative associations*: some diseases that are negatively associated with ASD with respect to normalized prevalence, *i.e.*, having those codes relatively over-represented in one's diagnostic history favors ending up in the control cohort, 2) *impact of sex*: there are sex-specific differences in the impact of specific disorders, and given a fixed level of impact, the number of codes that drive the outcomes is significantly more in males (Fig. 5A vs B).

Some of the disorders that show up in Fig. 5, panels A and B are surprising, *e.g.*, congenital hemiplegia or diplegia of the upper limbs indicative of either cerebral palsy (CP) or a spinal cord/brain injury, neither of which has a direct link to autism. Since only about 7% of the children with cerebral palsy (CP) are estimated to have a co-occurring ASD (47, 48), and with the prevalence of CP significantly lower (1 in 352 vs 1 in 59 for autism), it follows that only a small number of children (approximately 1.17%) with autism have co-occurring CP. Thus, with significantly higher prevalence in children diagnosed with autism compared to the general population (1.7% vs 0.28%), CP codes show up with higher odds in the true positive set. Also, SI-Fig. 1A shows that the immunological, metabolic, and endocrine disorders are almost completely risk-increasing. In contrast, respiratory diseases (panel B) are largely risk-decreasing. On the other hand, infectious diseases have roughly equal representations in the risk-increasing and risk-decreasing classes (panel C). The risk-decreasing infectious diseases tend to be due to viral or fungal organisms, which might point to the use of antibiotics in bacterial infections, and the consequent dysbiosis of the gut microbiota (23, 37) as a risk factor.

Any predictive analysis of ASD must address if we can discriminate ASD from general developmental and behavioral disorders. The DSM-5 established a single category of ASD to replace the subtypes of autistic disorder, Asperger syndrome, and pervasive developmental disorders (8). This aligns with our use of diagnostic codes from ICD9 299.X as specification of an ASD diagnosis, and use standardized mapping to 299.X from ICD10 codes when we encounter them. For other psychiatric disorders, we get high discrimination reaching AUCs over 90% at 100 – 125 weeks of age (SI-Fig. 5A), which establishes that our pipeline is indeed largely specific to ASD.

We carried out a battery of tests to ensure that our results are not significantly impacted by class imbalance (since our

178 more than 1% of cases of idiopathic autism (46).

179 In the absence of biomarkers, current screening in pediatric
180 primary care visits uses standardized questionnaires to cate-

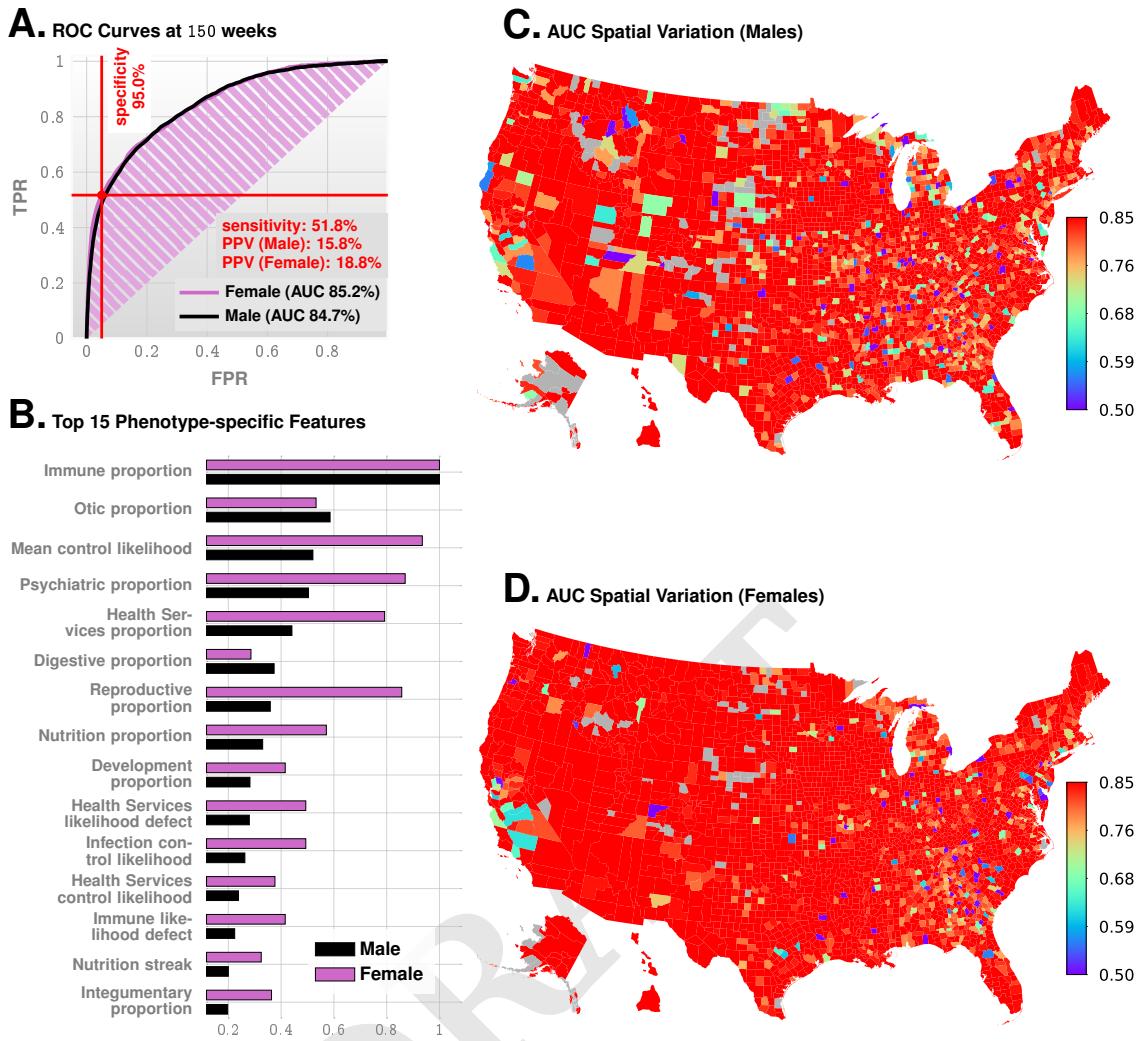


Fig. 2. Predictive Performance. Panel A shows the ROC curves for males and females. Panel B shows the feature importance inferred by our prediction pipeline. The detailed description of the features is given in Table 2. The most import feature is related to immunologic disorders, and we note that in addition to features related to individual disease categories, we also have the mean control likelihood (rank 3), which may be interpreted as the average likelihood of the diagnostic patterns correspond to the control category as opposed to the positive category. Panels C and D show the spatial variation in the achieved predictive performance at 150 weeks, measured by AUC, for males and females, respectively. Gray areas lack data on either positive or negative cases. These county-specific AUC plots show that the performance of the algorithm has relatively weak geospatial dependence, which is important in the light of current uneven distribution of diagnostic resources.

control cohort is orders of magnitude larger) or systematic coding errors (See Methods), *e.g.*, we verified that restricting the positive cohort to children with at least two distinct ASD diagnostic codes in their medical histories instead of one, has little impact on out-of-sample predictive performance (SI-Fig. 5B).

Can our performance be matched by simply asking how often a child is sick? We found that the density of codes in a child's medical history is indeed somewhat predictive of a future ASD diagnosis, with the $AUC \approx 75\%$ in the Truven database at 150 weeks (See SI-Fig. 5, panel D in the supplementary text). This is expected, since children with autism do indeed have higher rates of co-morbidities. However, it does not have stable performance across databases, and has no significant effect once the rest of the features are combined.

As a key limitation to our approach, automated pattern recognition might not reveal true causal precursors. The relatively uncurated nature of the data does not correct for coding mistakes by the clinician and other artifacts, *e.g.* a

bias towards over-diagnosis of children on the borderline of the diagnostic criteria due to clinicians' desire to help families access service, and biases arising from changes in diagnostic practices over time (49). Discontinuities in patient medical histories from change in provider-networks can also introduce uncertainties in risk estimates, and socio-economic status of patients which impact access to healthcare might skew patterns in EHR databases. Despite these limitations, the design of a questionnaire-free component to ASD screening that systematically leverages co-morbidities has far-reaching consequences, by potentially slashing the false positives and wait-times, as well as removing systemic under-diagnosis issues amongst females and minorities.

Future efforts will attempt to realize our approach within a clinical setting. We will also explore the impact of maternal medical history, and the use of calculated risk to trigger blood-work to look for expected transcriptomic signatures of ASD. Finally, the analysis developed here applies to phenotypes beyond ASD, thus opening the door to the possibility

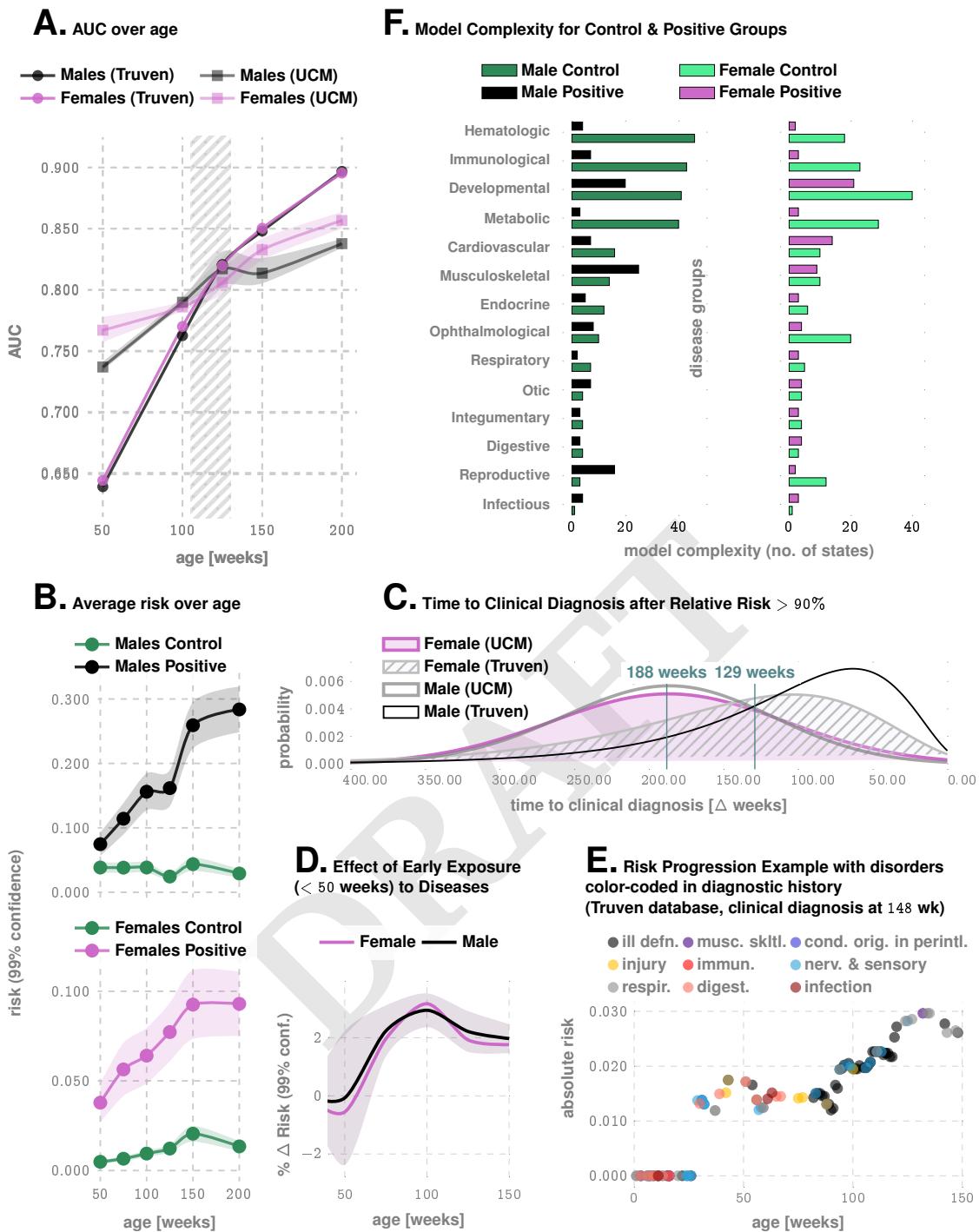


Fig. 3. More details on Predictive Performance and Variation of Inferred Risk. Panel A illustrates AUC achieved as a function of patient age, for the Truven and UCM datasets. The shaded area outlines the 2 - 2.5 years of age, and shows that we achieve > 80% AUC for either sex from shortly after 2 years. Panel B illustrates how the average risk changes with time for the control and the positive cohorts. Panel C shows the distribution of the prediction horizon: the time to a clinical diagnosis after inferred relative risk crosses 90%. Panel D shows that for each new disease code for a low-risk child, ASD risk increases by approximately 2% for either sex. Panel E illustrates the risk progression of a specific, ultimately autistic male child in the Truven database. Abbreviations in the legend: ill. defn. (Symptoms, Signs, And Ill-Defined Conditions), musc. skltl. (Diseases Of The Musculoskeletal System And Connective Tissue), cond. orig. in perintl. (Certain Conditions Originating In The Perinatal Period), immun. (Endocrine, Nutritional And Metabolic Diseases, And Immunity Disorders), nerv. & sensory (Diseases Of The Nervous System And Sense Organs), respir. (Respiratory Disorders), and digest. (Digestive Disorders). Panel F illustrates how inferred models differ between the control vs. the positive cohorts. On average, models get less complex, implying the exposures get more statistically independent.

284 of general comorbidity-aware risk predictions from electronic
285 health record databases.

Data Sharing. Software implementation of the pipeline is available at: <https://github.com/zeroknowledgediscovery/ehrzero>, and installation in standard python environments may be

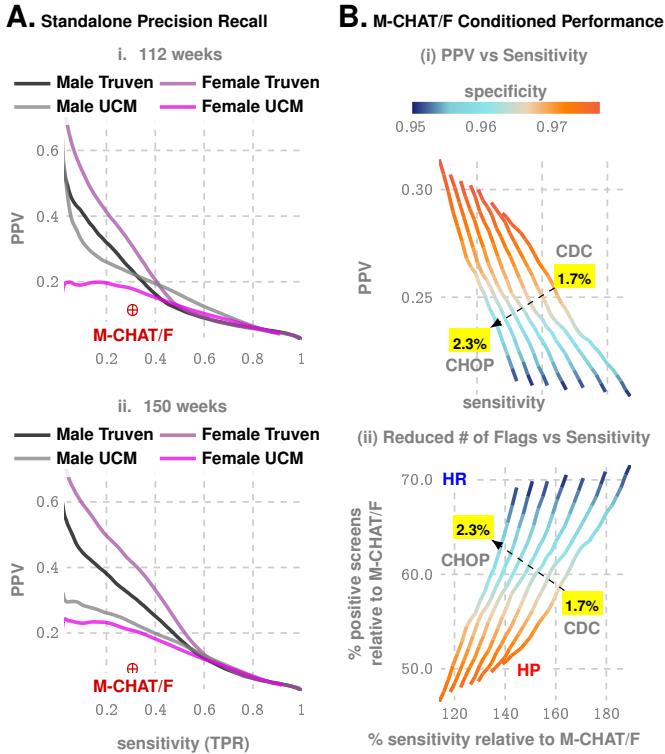


Fig. 4. Metrics relevant to clinical practice: PPV vs Sensitivity trade-offs. Panel A shows the precision/recall curves, *i.e.*, the trade-off between PPV and sensitivity. Panel B shows how we can boost performance using population stratification from the distribution of M-CHAT/F scores in the population, as reported by the CHOP study (3). Panel C illustrates the boosted performance compared to M-CHAT/F alone, measured by the relative percentage increase in sensitivity, and percentage decrease in positive screens. Note that the population prevalence impacts this optimization, and hence we have a distinct curve for each prevalence value (1.7% is the CDC estimate, while 2.23% is reported by the CHOP study). The two extreme operating zones marked as High Precision (HP) and High Recall (HR): if we choose to operate in HR, then we do not reduce the number of positive screens by much, but maximize sensitivity, while by operating in HP, we do not increase sensitivity by much but double the PPV achieved in current practice. Note in all these zones, we maintain specificity above 95%, which is the current state of art, implying that by doubling the PPV, we can halve the number of positive screens currently reported, thus potentially sharply reducing the queues and wait-times.

done from <https://pypi.org/project/ehrzero/>. A sample of de-identified data from the UCM database is be shared as part of the public software package.

Materials and Methods

Source of Electronic Patient Records. We view the task of predicting ASD diagnosis as a binary classification problem: sequences of diagnostic codes are classified into positive and control categories, where “positive” refers to children eventually diagnosed with ASD, as indicated by the presence of a clinical diagnosis (ICD9 code 299.X) in their medical records. Of the two independent sources of clinical incidence data used in this study, the primary source used to train our predictive pipeline is the Truven Health Analytics MarketScan® Commercial Claims and Encounters Database for the years 2003 to 2012 (50) (referred to as the Truven dataset). This US national database contains data contributed by over 150 insurance carriers, large self-insuring companies, and Medi and is a culmination of over 4.6 billion inpatient and outpatient service claims and almost six billion diagnosis codes. We extracted histories of patients within the age of 0 – 9 years, and excluded patients for whom: 1) At least one code of any available phenotypes is present, 2) Lag between first and last available record for a patient should

be at least 15 weeks. These exclusion criteria ensure that we are not considering patients who have too few observations to either train on. Additionally, during validation runs, we restricted the control set to patients observable in the databases to those whose last record is not before the first 150 weeks of life. Characteristics of excluded patients is shown in Table 1. We trained with over 30M diagnostic records (16,649,548 for males and 14,318,303 for females with 9,835 unique codes).

While the Truven database is used for both training and out-of-sample cross-validation with held-back data, our second independent dataset consisting of de-identified diagnostic records for children treated at the University of Chicago Medical Center between the years of 2006 to 2018 (the UCM dataset), aids in further cross-validation. We considered children between the ages of 0 – 5 years, and applied the same exclusion criteria as the Truven dataset. The number of patients used from the two databases is shown in Table 1.

Our datasets are consistent with documented prevalence, with no significant geospatial prevalence variation (Fig. 1D) and reveal that infections and immunological disorders have differential representation in the positive and control groups (Fig. 1C). The median diagnosis age is just over 3 years in the claims database (Fig. 1B) versus 3 years 10 months to 4 years in US (51). Cohort details are given in Table 1 and discussed in Methods. Importantly, for the positive cohort, we only consider diagnostic history up to the first ASD code.

The significant diversity of diagnostic codes (6,089 distinct ICD-9 codes and 11,522 distinct ICD-10 codes in total in the two datasets), along with the sparsity of codes per sequence and the need to make good predictions as early as possible, makes this a difficult learning problem, where standard deep learning approaches do not suffice (See Table 2). To address these issues, we proceed by partitioning the disease spectrum into 17 broad categories, *e.g.* infectious diseases, immunologic disorders, endocrinial disorders etc. Each patient is then represented by 17 distinct time series, each tracking an individual disease category. At the population level, these disease-specific sparse stochastic time series are compressed into specialized Markov models (separately for the control and the treatment cohorts) to identify the distinctive patterns pertaining to elevated ASD risk. With these inferred patterns included as features (Table 3) we train a second level predictor that learns to map individual patients to the control or the positive groups based on their similarity to the identified Markov models of category-specific diagnostic histories (See Methods).

Time-series Modeling of Diagnostic History. Individual diagnostic histories can have long-term memory (52), implying that the order, frequency, and comorbid interactions between diseases are important for assessing the future risk of our target phenotype. We analyze patient-specific diagnostic code sequences by first representing the medical history of each patient as a set of stochastic categorical time-series | one each for a specific group of related disorders | followed by the inference of stochastic models for these individual data streams. These inferred generators are from a special class of Hidden Markov Models (HMMs), referred to as Probabilistic Finite State Automata (PFSA) (53). The inference algorithm we use is distinct from classical HMM learning, and has important advantages related to its ability to infer structure, and its sample complexity (See Supplementary text, Section 11). We infer a separate class of models for the positive and control cohorts, and then the problem reduces to determining the probability that the short diagnostic history from a new patient arises from the positive as opposed to the control category of the inferred models.

Step 1: Partitioning The Human Disease Spectrum. We begin by partitioning the human disease spectrum into 17 non-overlapping categories, as shown in Table 2. Each category is defined by a set of diagnostic codes from the International Classification of Diseases, Ninth Revision (ICD9) (See Table 2 in the main text and Table SI-4 in the Supplementary text for description of the categories used in this study). For this study, we considered 9,835 distinct ICD9 codes (and their ICD10 General Equivalence Mappings (GEMS) (17) equivalents). We came across 6,089 distinct ICD-9 codes and 11,522 distinct ICD-10 codes in total in the two datasets we analyzed. 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;

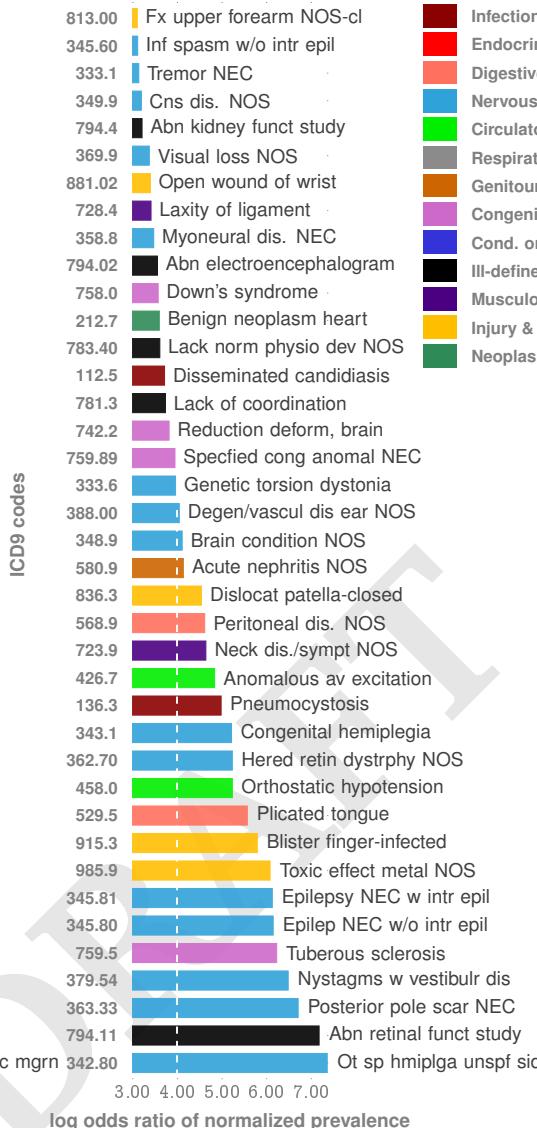
A. Male (3 YR)

782.0	Skin sensation disturb
379.54	Nystagms w vestibul dis
784.61	Alexia and dyslexia
366.50	After-cataract NOS
458.0	Orthostatic hypotension
743.32	Cortical/zonular catarac
906.5	Late eff head/neck burn
779.9	Perinatal condition NOS
345.60	Inf spasm w/o intr epil
779.8	Neonatal bradycardia
783.4	Lack norm physio dev NOS
264.9	Vitamin A deficiency NOS
874.3	Open wound thyroid-compl
722.6	Disc degeneration NOS
378.63	Mech strab w oth conditn
387.2	Cochlear otosclerosis
756.4	Chondrodyostrophy
784.60	Symbolic dysfunction NOS
345.81	Epilepsy NEC w intr epil
345.61	Inf spasm w intract epil
078.5	Cytomegaloviral disease
794.14	Abn oculomotor studies
715.95	Osteoarthros NOS-pelvis
374.03	Spastic entropion
337.9	Autonomic nerve dis NEC
343.4	Infantile hemiplegia
362.76	Vitelliform dystrophy
191.7	Mal neo brain stem
344.2	Diplegia of upper limbs
696.0	Psoriatic arthropathy
271.9	Dis carbohydr metab NOS
230.3	Ca in situ colon
556.1	Ulcerative ileocolitis
227.3	Benign neo pituitary
345.71	Epil par cont w intr epi
727.68	Rupture tendon foot NEC
985.9	Toxic effect metal NOS
364.02	Recurrent iridocyclitis
531.40	Chr stomach ulc w hem
346.80	Othr migrne wo ntrc mgnr

4.00 5.00 6.00 7.00

log odds ratio of normalized prevalence

B. Female (3 YR)

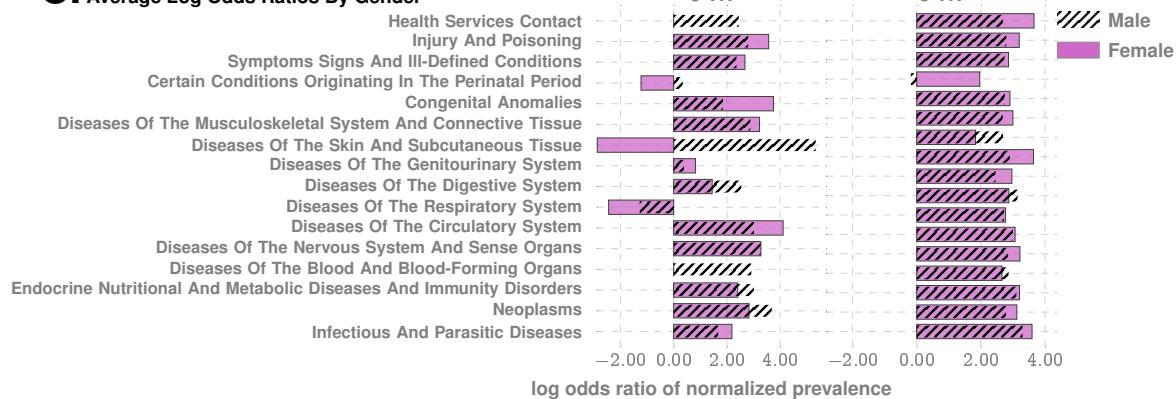


3.00 4.00 5.00 6.00 7.00
log odds ratio of normalized prevalence

ICD9 Class



C. Average Log Odds Ratios By Gender



-2.00 0.00 2.00 4.00 -2.00 0.00 2.00 4.00
log odds ratio of normalized prevalence

Fig. 5. Co-morbidity Patterns Panel A and B. Difference in occurrence frequencies of diagnostic codes between true positive (TP) and true negative (TN) predictions. The dotted line on panel B shows the abscissa lower cut-off in Panel A, illustrating the lower prevalence of codes in females. Panel C illustrates log-odds ratios for ICD9 disease categories at different ages. Importantly, the negative associations disappear when we consider older children, consistent with the lack of such reports in the literature which lack studies on very young cohorts.

388 we want our algorithm to seek out the important patterns without
389 any manual curation of the input data. The limitation of the set of
390 phenotypes to 9835 unique codes arises from excluding patients from
391 the database who have very few and rare codes that will skew the
392 statistical estimates. As shown in Table 1, we exclude a very small
393 number of patients, and who have very short diagnostic histories
394 with a very small number of codes.

For each patient, the past medical history is a sequence $(t_1, x_1), \dots, (t_m, x_m)$, 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:

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

395 The time-series z^k is terminated at a particular week if the patient
396 is diagnosed with ASD the week after. Thus for patients in the
397 control cohort, the length of the mapped trinary series is limited by
398 the time for which the individual is observed within the 2003 – 2012
399 span of our database. In contrast, for patients in the positive cohort,
400 the length of the mapped series reflect the time to the first ASD
401 diagnosis. Patients do not necessarily enter the database at birth,
402 and we prefix each series with 0s to approximately synchronize
403 observations to age in weeks. Each patient is now represented by
404 17 mapped trinary series.

Step 2: Model Inference & The Sequence Likelihood Defect. The
405 mapped series, stratified by sex, disease-category, and ASD
406 diagnosis-status are considered to be independent sample paths,
407 and we want to explicitly model these systems as specialized HMMs
408 (PFSA). We model the positive and the control cohorts for each
409 sex, and in each disease category separately, ending up with a total
410 of 68 HMMs at the population level (17 categories, 2 sexes, 2
411 cohort-types: positive and control, SI-Fig. 7 in the supplementary
412 text provides some examples). Each of these inferred models is a
413 PFSA; a directed graph with probability-weighted edges, and
414 acts as an optimal generator of the stochastic process driving the
415 sequential appearance of the three letters (as defined by Eq. Eq. (1))
416 corresponding to each sex, disease category, and cohort-type (See
417 Section 11 in the Supplementary text for background on PFSA
418 inference).

To reliably infer the cohort-type of a new patient, *i.e.*, the likelihood of a diagnostic sequence being generated by the corresponding cohort model, we generalize the notion of Kullbeck-Leibler (KL) divergence (54, 55) between probability distributions to a divergence $\mathcal{D}_{\text{KL}}(G||H)$ between ergodic stationary categorical stochastic processes (56) G, H as:

$$\mathcal{D}_{\text{KL}}(G||H) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x:|x|=n} p_G(x) \log \frac{p_G(x)}{p_H(x)} \quad [2]$$

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 :

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

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 Suppl. text Section 11, Theorem 6 and 7):

$$\lim_{|x| \rightarrow \infty} L(x, G) = \mathcal{H}(G) \quad [4a]$$

$$\lim_{|x| \rightarrow \infty} L(x, H) = \mathcal{H}(H) + \mathcal{D}_{\text{KL}}(G||H) \quad [4b]$$

where $\mathcal{H}(\cdot)$ is the entropy rate of a process (54). Importantly, Eq. 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 ASD, 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, Δ^j) as:

$$\Delta^j \triangleq L(G_0^j, x) - L(G_+^j, x) \rightarrow \mathcal{D}_{\text{KL}}(G_0^j || G_+^j) \quad [5]$$

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

Step 3: Risk Estimation Pipeline With Semi-supervised & Supervised Learning Modules. The risk estimation pipeline operates on patient specific information limited to the sex and available diagnostic history from birth, and produces an estimate of the relative risk of ASD diagnosis at a specific age, 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 (57). The complete list of 165 features used is provided in Table 3.

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. 2B shows the top 15 features ranked in order of their relative importance (relative loss in performance when dropped out of the analysis).

Calculating Relative Risk. Our pipeline maps medical histories to a raw indicator of risk. However, to make crisp predictions, we must choose a decision threshold for this raw score. Conceptually identical to the notion of Type 1 and Type 2 errors in classical statistical analyses, the choice of a threshold trades off false positives (Type 1 error) for false negatives (Type 2 error): choosing a small threshold results in predicting a larger fraction of future diagnoses correctly, *i.e.* have a high true positive rate (TPR), while simultaneously suffering from a higher false positive rate (FPR), and vice versa. Therefore, a choice of a specific decision threshold reflects a choice of the maximum FPR and minimum TPR, and is driven by the application at hand. In this study, we base our analysis on maximizing the F_1 -score, defined as the harmonic mean of sensitivity and specificity, to make a balanced trade-off between the two kinds of errors (See Supplementary text, Section 5). The *relative risk* is then defined as the ratio of the raw risk to the decision threshold, and a value > 1 predicts a future ASD diagnosis.

Boosting Performance Via Leveraging Population Stratification Induced By Existing Tests. We leverage the population stratification induced by an existing independent screening test (M-CHAT/F) to improve combined performance. Here a combination refers to the conditional choice of the sensitivity/specificity trade-offs for our tool in each sub-population such that the overall performance is optimized with respect to whether we wish to maximize the PPV or the sensitivity at a specified minimum level of specificity. Assume that there are m sub-populations such that: the sensitivities, specificities achieved, and the prevalences in each sub-population are given by s_i, c_i and ρ_i respectively, with $i \in \{1, \dots, m\}$. Let β_i be the relative size of each sub-population. Then, we have (See Supplementary text, Section 9.1):

$$s = \sum_{i=1}^m s_i \gamma_i \quad [6a]$$

$$c = \sum_{i=1}^m c_i \gamma'_i \quad [6b]$$

where we have denoted:

$$\gamma_i = \beta_i \frac{\rho_i}{\rho}, \text{ and } \gamma'_i = \beta_i \frac{1 - \rho_i}{1 - \rho} \quad [6c]$$

and s, c, ρ are the overall sensitivity, specificity, and prevalence. Knowing the values of γ_i, γ'_i , we can carry out an m -dimensional search to identify the feasible choices of s_i, c_i pairs for each i , such that some global constraint is satisfied, *e.g.* minimum values of

specificity, sensitivity, and PPV. We consider 4 sub-populations defined by M-CHAT/F score brackets (3), and if the screen result is considered a positive (high risk, indicating the need for a full diagnostic evaluation) or a negative, *i.e.*, low risk: 1) score ≤ 2 screening ASD negative, 2) score [3 – 7] screening ASD negative on follow-up, 3) score [3 – 7] and screening ASD positive on follow-up, and 4) score ≥ 8 , screening ASD positive. (See SI-Table 2). The “follow-up” in the context of M-CHAT/F refers to the re-evaluation of responses by qualified personnel. We use published data on the relative sizes and the prevalence statistics in these sub-populations (3) to compute the feasible conditional choices of our operating point to strictly supersede M-CHAT/F performance. Two limiting operating conditions are of special interest here, where we maximize PPV under some minimum specificity and sensitivity (denoted as the High Precision or the HP operating point), and where we maximize sensitivity under some minimum PPV and specificity (denoted as the High Recall or the HR operating point). Taking these minimum values of specificity, sensitivity, and PPV to be those reported for M-CHAT/F, we identify the set feasible set of conditional choices in a four-dimensional decision space that would outperform M-CHAT/F in universal screening. The results are shown in Fig. 4B.

ACKNOWLEDGMENTS. This work is funded in part by the Defense Advanced Research Projects Agency (DARPA) project #FP070943-01-PR. The claims made in this study do not reflect the position or the policy of the US Government. The UCM dataset is provided by the Clinical Research Data Warehouse (CRDW) maintained by the Center for Research Informatics (CRI) at the University of Chicago. The Center for Research Informatics is funded by the Biological Sciences Division, the Institute for Translational Medicine/CTSA (NIH UL1 TR000430) at the University of Chicago.

1. Data & statistics on autism spectrum disorder | cdc (2019).
2. L Gilotty, Early screening for autism spectrum (2019).
3. W Guthrie, et al., Accuracy of Autism Screening in a Large Pediatric Network. *Pediatrics* **144** (2019).
4. E Gordon-Lipkin, J Foster, G Peacock, Whittling Down the Wait Time: Exploring Models to Minimize the Delay from Initial Concern to Diagnosis and Treatment of Autism Spectrum Disorder. *Pediatr. Clin. North Am.* **63**, 851–859 (2016).
5. Data & statistics on autism spectrum disorder | cdc (2019).
6. LA Schieve, et al., Population attributable fractions for three perinatal risk factors for autism spectrum disorders, 2002 and 2008 autism and developmental disabilities monitoring network. *Ann Epidemiol* **24**, 260–266 (2014).
7. F Volkmar, et al., Practice parameter for the assessment and treatment of children and adolescents with autism spectrum disorder. *J. Am. Acad. Child & Adolesc. Psychiatry* **53**, 237–257 (2014).
8. SL Hyman, SE Levy, SM Myers, , et al., Identification, evaluation, and management of children with autism spectrum disorder. *Pediatrics* **145** (2020).
9. LG Kalb, et al., Determinants of appointment absenteeism at an outpatient pediatric autism clinic. *J. Dev. & Behav. Pediatr.* **33**, 685–697 (2012).
10. J Bisgaier, D Levinson, DB Cutts, KV Rhodes, Access to autism evaluation appointments with developmental-behavioral and neurodevelopmental subspecialists. *Arch. pediatrics & adolescent medicine* **165**, 673–674 (2011).
11. TS Fenikilé, K Ellerbeck, MK Filippi, CM Daley, Barriers to autism screening in family medicine practice: a qualitative study. *Prim. health care research & development* **16**, 356–366 (2015).
12. E Gordon-Lipkin, J Foster, G Peacock, Whittling down the wait time: exploring models to minimize the delay from initial concern to diagnosis and treatment of autism spectrum disorder. *Pediatr. Clin. North Am.* **63**, 851–859 (2016).
13. DL Robins, et al., Validation of the modified checklist for autism in toddlers, revised with follow-up (m-chat-rf). *Pediatrics* **133**, 37–45 (2014).
14. C Tye, AK Runicles, AJO Whitehouse, GA Alvares, Characterizing the Interplay Between Autism Spectrum Disorder and Comorbid Medical Conditions: An Integrative Review. *Front Psychiatry* **9**, 751 (2018).
15. IS Kohane, et al., The co-morbidity burden of children and young adults with autism spectrum disorders. *PLoS ONE* **7**, e33224 (2012).
16. VG Jarquin, LD Wiggins, LA Schieve, K Van Naarden-Braun, Racial disparities in community identification of autism spectrum disorders over time; metropolitan atlanta, georgia, 2000–2006. *J. Dev. & Behav. Pediatr.* **32**, 179–187 (2011).
17. General equivalence mappings (year?).
18. LA Althouse, JA Stockman, Pediatric workforce: A look at pediatric nephrology data from the american board of pediatrics. *The J. pediatrics* **148**, 575–576 (2006).
19. O Zerbo, et al., Immune mediated conditions in autism spectrum disorders. *Brain Behav. Immun.* **46**, 232–236 (2015).
20. H Won, W Mah, E Kim, Autism spectrum disorder causes, mechanisms, and treatments: focus on neuronal synapses. *Front Mol Neurosci* **6**, 19 (2013).
21. G Xu, et al., Association of Food Allergy and Other Allergic Conditions With Autism Spectrum Disorder in Children. *JAMA Netw Open* **1**, e180279 (2018).

22. JB Adams, et al., Nutritional and metabolic status of children with autism vs. neurotypical children, and the association with autism severity. *Nutr Metab (Lond)* **8**, 34 (2011).
23. A Fattorusso, L Di Genova, GB Dell’Isola, E Mencaroni, S Esposito, Autism Spectrum Disorders and the Gut Microbiota. *Nutrients* **11** (2019).
24. R Diaz Heijtz, et al., Normal gut microbiota modulates brain development and behavior. *Proc. Natl. Acad. Sci. U.S.A.* **108**, 3047–3052 (2011).
25. S Rose, et al., Mitochondrial dysfunction in the gastrointestinal mucosa of children with autism: A blinded case-control study. *PLoS ONE* **12**, e0186377 (2017).
26. EM Sajdel-Sulkowska, et al., Common Genetic Variants Link the Abnormalities in the Gut-Brain Axis in Prematurity and Autism. *Cerebellum* **18**, 255–265 (2019).
27. MS Kayser, J Dalmat, Anti-NMDA Receptor Encephalitis in Psychiatry. *Curr Psychiatry Rev* **7**, 189–193 (2011).
28. OI Dadalko, BG Travers, Evidence for Brainstem Contributions to Autism Spectrum Disorders. *Front Integr Neurosci* **12**, 47 (2018).
29. Y Yamashita, et al., Anti-Inflammatory Effect of Ghrelin in Lymphoblastoid Cell Lines From Children With Autism Spectrum Disorder. *Front Psychiatry* **10**, 152 (2019).
30. L Shen, et al., Proteomics Study of Peripheral Blood Mononuclear Cells (PBMCs) in Autistic Children. *Front Cell Neurosci* **13**, 105 (2019).
31. K Ohja, et al., Neuroimmunologic and Neurotrophic Interactions in Autism Spectrum Disorders: Relationship to Neuroinflammation. *Neuromolecular Med.* **20**, 161–173 (2018).
32. D Gadyasz, A Krzywdziska, KK Hozyasz, Immune Abnormalities in Autism Spectrum Disorder—Could They Hold Promise for Causative Treatment? *Mol. Neurobiol.* **55**, 6387–6435 (2018).
33. TC Theoharides, I Tsilioni, AB Patel, R Doyle, Atopic diseases and inflammation of the brain in the pathogenesis of autism spectrum disorders. *Transl Psychiatry* **6**, e844 (2016).
34. AM Young, et al., From molecules to neural morphology: understanding neuroinflammation in autism spectrum condition. *Mol Autism* **7**, 9 (2016).
35. LA Croen, et al., Family history of immune conditions and autism spectrum and developmental disorders: Findings from the study to explore early development. *Autism Res* **12**, 123–135 (2019).
36. T Vargason, DL McGuinness, J Hahn, Gastrointestinal Symptoms and Oral Antibiotic Use in Children with Autism Spectrum Disorder: Retrospective Analysis of a Privately Insured U.S. Population. *J Autism Dev Disord* **49**, 647–659 (2019).
37. M Fiorentino, et al., Blood-brain barrier and intestinal epithelial barrier alterations in autism spectrum disorders. *Mol Autism* **7**, 49 (2016).
38. FK Satterstrom, et al., Large-scale exome sequencing study implicates both developmental and functional changes in the neurobiology of autism. *bioRxiv* (2019).
39. T Gaugler, et al., Most genetic risk for autism resides with common variation. *Nat. Genet.* **46**, 881–885 (2014).
40. DL Vargas, C Nasimbeni, C Krishnan, AW Zimmerman, CA Pardo, Neuroglial activation and neuroinflammation in the brain of patients with autism. *Ann. Neurol.* **57**, 67–81 (2005).
41. H Wei, et al., IL-6 is increased in the cerebellum of autistic brain and alters neural cell adhesion, migration and synaptic formation. *J Neuroinflammation* **8**, 52 (2011).
42. AM Young, E Campbell, S Lynch, J Suckling, SJ Powis, Aberrant NF-kappaB expression in autism spectrum condition: a mechanism for neuroinflammation. *Front Psychiatry* **2**, 27 (2011).
43. HK Hughes, E Mills Ko, D Rose, P Ashwood, Immune Dysfunction and Autoimmunity as Pathological Mechanisms in Autism Spectrum Disorders. *Front Cell Neurosci* **12**, 405 (2018).
44. N Pearce, The ecological fallacy strikes back. *J. epidemiology community health* **54**, 326–7 (2000).
45. JD Murdoch, MW State, Recent developments in the genetics of autism spectrum disorders. *Curr. Opin. Genet. Dev.* **23**, 310–315 (2013).
46. VW Hu, The expanding genomic landscape of autism: discovering the ‘forest’ beyond the ‘trees’. *Futur. Neural* **8**, 29–42 (2013).
47. Prevalence of cerebral palsy, co-occurring autism spectrum disorders, and motor functioning (2020).
48. D Christensen, et al., Prevalence of cerebral palsy, co-occurring autism spectrum disorders, and motor functioning—a us and d developmental disabilities monitoring n etwork, usa, 2008. *Dev. Medicine & Child Neurol.* **56**, 59–65 (2014).
49. EM Rødgaard, K Jensen, JN Vergnes, I Soulières, L Mottron, Temporal Changes in Effect Sizes of Studies Comparing Individuals With and Without Autism: A Meta-analysis. *JAMA Psychiatry* **76**, 1124–1132 (2019).
50. L Hansen, The truven health marketscan databases for life sciences researchers. *Truven Heal. Analytics IBM Watson Heal.* (2017).
51. J Bai, et al., Prevalence of Autism Spectrum Disorder Among Children Aged 8 Years - Autism and Developmental Disabilities Monitoring Network, 11 Sites, United States, 2014. *MMWR Surveill Summ* **67**, 1–23 (2018).
52. CWJ Granger, R Joyeux, An introduction to long-memory time series models and fractional differencing. *J. Time Ser. Analysis* **1**, 15–29 (year?).
53. I Chattopadhyay, H Lipson, Abductive learning of quantized stochastic processes with probabilistic finite automata. *Philos Trans A* **371**, 20110543 (2013).
54. TM Cover, JA Thomas, *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. (Wiley-Interscience, New York, NY, USA), (2006).
55. S Kulback, RA Leibler, On information and sufficiency. *Ann. Math. Stat.* **22**, 79–86 (1951).
56. J Doob, *Stochastic Processes*, Wiley Publications in Statistics. (John Wiley & Sons), (1953).
57. JH Friedman, Stochastic gradient boosting. *Comput. Stat. Data Anal.* **38**, 367–378 (2002).

Supplementary Information:

Reduced False Positives in Autism Screening Via Digital Bio-markers Inferred from Deep Co-morbidity Patterns

Dmytro Onishchenko^a, Yi Huang^a, James van Horne^a, Peter J. Smith^{d,g}, Michael M. Msall^{e,f}, Ishanu Chattopadhyay^{a,b,c,*}

^a*Department of Medicine, University of Chicago, Chicago, IL, USA*

^b*Committee on Genetics, Genomics & Systems Biology, University of Chicago, Chicago, IL, USA*

^c*Committee on Quantitative Methods in Social, Behavioral, and Health Sciences, University of Chicago, Chicago, IL, USA*

^d*Department of Pediatrics, Section of Developmental and Behavioral Pediatrics, University of Chicago, Chicago, IL, USA*

^e*Department of Pediatrics, Section Chief of Developmental and Behavioral Pediatrics, University of Chicago, Chicago, IL, USA*

^f*Joseph P. Kennedy Research Center on Intellectual and Neurodevelopmental Disabilities, University of Chicago, Chicago, IL, USA*

^g*Executive Committee Chair, American Academy of Pediatrics' Section on Developmental and Behavioral Pediatrics*

Contents

1 Pipeline Optimization	16
1.1 Input Data Format	16
1.2 Algorithms	16
2 Example Run with Released Application	19
2.1 Prerequisites & Installation	19
2.2 EHR data format	20
2.3 Sample Python code risk estimation	20
2.4 Sample Python script risk estimation	20
3 Comparison With State of the Art Off-the-shelf ML Algorithms	20
4 Comparison With Pipeline Variations, Feature Subsets and Neural Net Post-processing	21
4.1 Feature Subset Evaluations & Code Density As A Feature	21
5 Threshold Selection on ROC Curve	21
6 Note on Reciever Operating Characteristics (ROC) and Precision-recall Curves	22
7 Effect of Class Imbalance	23
8 Note on ASD Clinical Diagnosis & Uncertainty of EHR Record	23
8.1 Diagnostic Evaluations	23

*To whom correspondence should be addressed: e-mail: ishanu@uchicago.edu.

8.2 Change In Diagnostic Criteria for ASD, Inclusion of PDD, Asperger, and Disambiguation From Unrelated Psychiatric Phenotypes	24
8.3 Performance Comparison With M-CHAT/F	24
9 Improving Wait-times For Diagnostic Evaluations by Reducing False Positives in Routine Screening	25
9.1 4D Decision Optimization Using M-CHAT/F Population Stratification To Boost PPV	25
10 Generating PFSA Models From Set of Input Streams with Variable Input Lengths	26
11 Probabilistic Finite State Automata Inference	26
11.1 Probabilistic Finite-State Automaton	26
12 Sequence Likelihood Defect	28

SI-Table 1: Boosted Sensitivity, Specificity and PPV Achieved at 150 weeks Conditioned on M-CHAT/F Scores

M-CHAT/F Outcome				global performance (Truven)			global performance (UCM)			prevalence
0-2 NEG	3-7 NEG	3-7 POS	≥ 8 POS	specificity	sensitivity	PPV	specificity	sensitivity	PPV	
specificity choices										
0.28	0.66	0.93	0.97	0.95	0.64	0.224	0.95	0.577	0.206	0.022
0.31	0.67	0.9	0.97	0.95	0.641	0.223	0.95	0.577	0.205	0.022
0.54	0.86	0.97	0.99	0.98	0.494	0.361	0.98	0.393	0.31	0.022
0.41	0.89	0.96	0.99	0.98	0.493	0.362	0.98	0.391	0.311	0.022
0.31	0.61	0.86	0.98	0.95	0.808	0.219	0.95	0.713	0.198	0.017
0.33	0.6	0.86	0.98	0.95	0.809	0.218	0.95	0.715	0.197	0.017
0.66	0.95	0.98	0.99	0.98	0.574	0.337	0.98	0.417	0.269	0.017
0.53	0.97	0.98	0.99	0.98	0.573	0.337	0.98	0.412	0.267	0.017
0.54	0.91	0.97	0.99	0.978	0.615	0.322	0.978	0.499	0.278	0.017
0.52	0.92	0.97	0.99	0.978	0.612	0.324	0.978	0.492	0.278	0.017

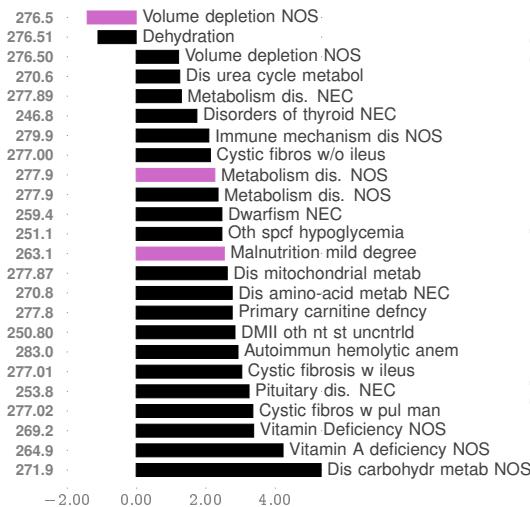
SI-Table 2: Population Stratification Results on large M-CHAT/F Study(n=20,375)¹

Id	Sub-population	Test Result	ASD positive	ASD Negative	Total %
A	M-CHAT/F ≥ 8	Positive	0.34%	0.64%	0.99%
B	M-CHAT/F $\in [3, 7]$	Positive (after follow-up)	0.52%	4.39%	4.91%
C	M-CHAT/F $\in [3, 7]$	Negative (after follow-up)	0.14%	3.1%	3.24%
D	M-CHAT/F $\in [0, 2]$	Negative	1.22%	89.63%	90.86%
Total %			2.23%	97.77%	100%

SI-Table 3: γ, γ' Computed from Population Stratification Recorded In M-CHAT/F Study¹ ($\rho = 0.0223$)

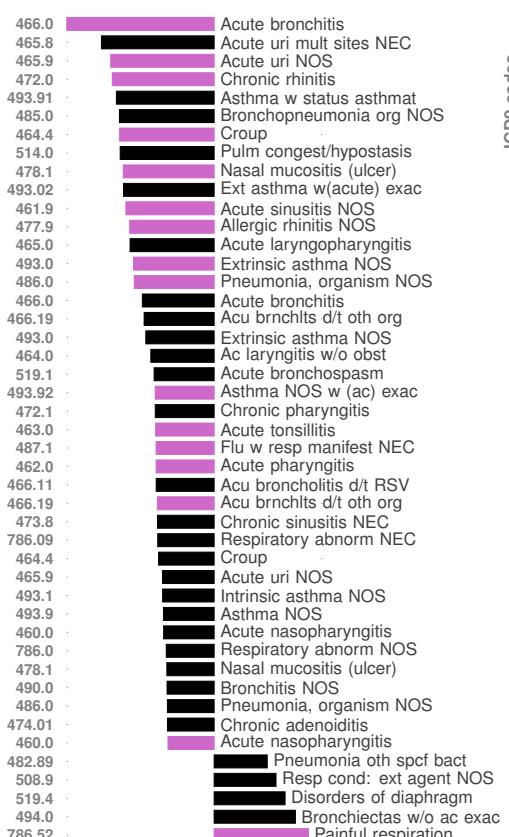
Id	Sub-population	Test Result	β_i	ρ_i	γ_i	γ'_i
A	M-CHAT/F ≥ 8	Positive	.0099	.3469	.1540	.0066
B	M-CHAT/F $\in [3, 7]$	Positive (after follow-up)	.0491	.1059	.2331	.0449
C	M-CHAT/F $\in [3, 7]$	Negative (after follow-up)	.0324	.0432	.0627	.0317
D	M-CHAT/F $\in [0, 2]$	Negative	.9086	.0134	.5471	.9168

A. Endocrine Nutritional Metabolic And Immunity Dis.



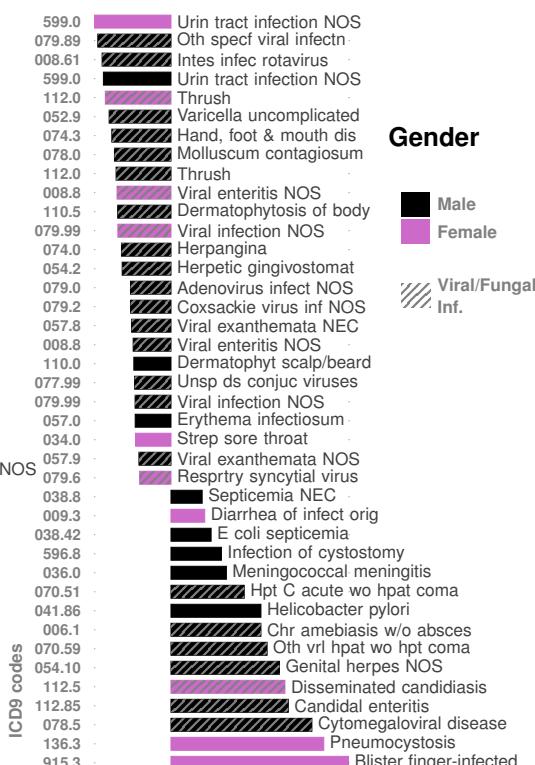
log odds ratio of normalized prevalence

B. Respiratory Disorders



log odds ratio of normalized prevalence

C. Infectious And Parasitic Diseases



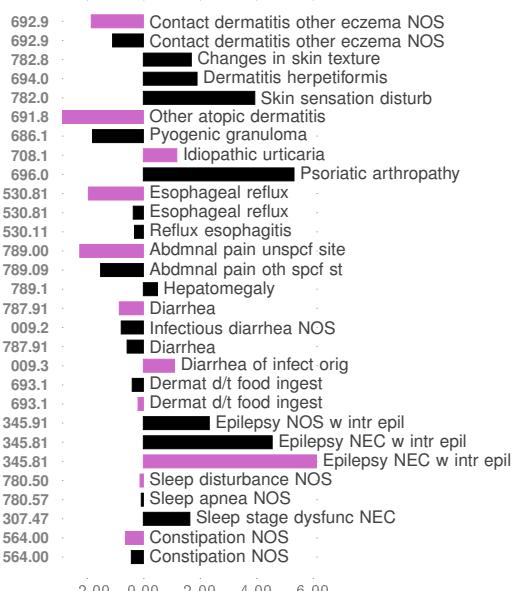
log odds ratio of normalized prevalence

Gender



Viral/Fungal Inf.

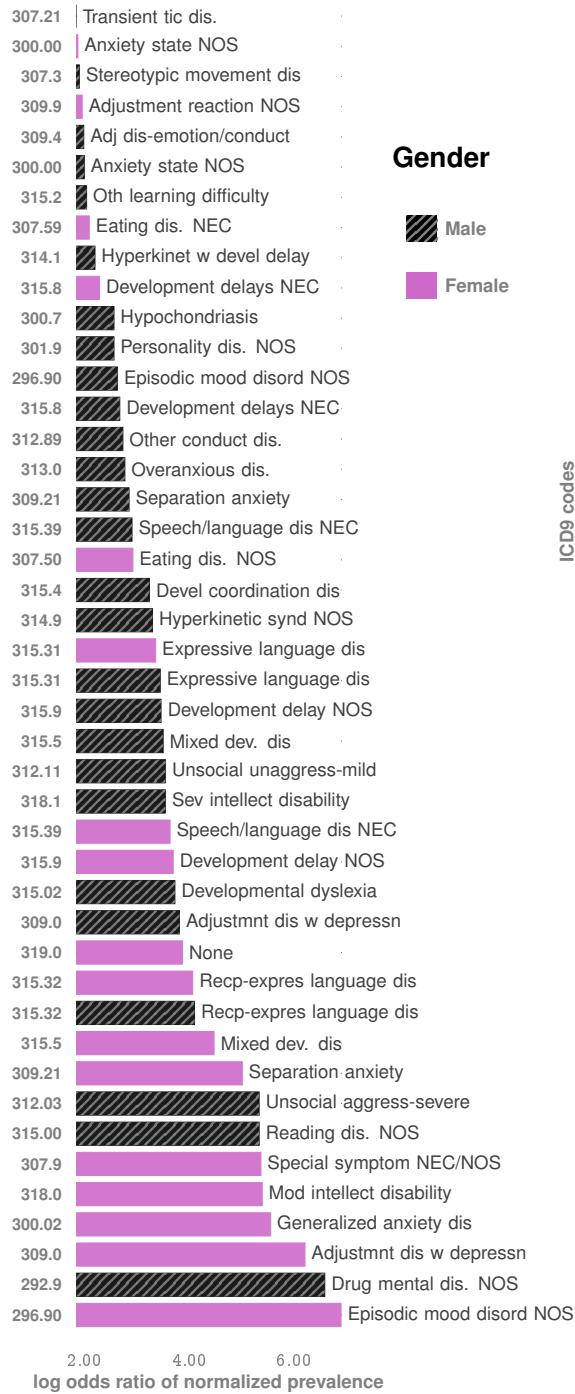
D. Similar Dis. with Opposed Association



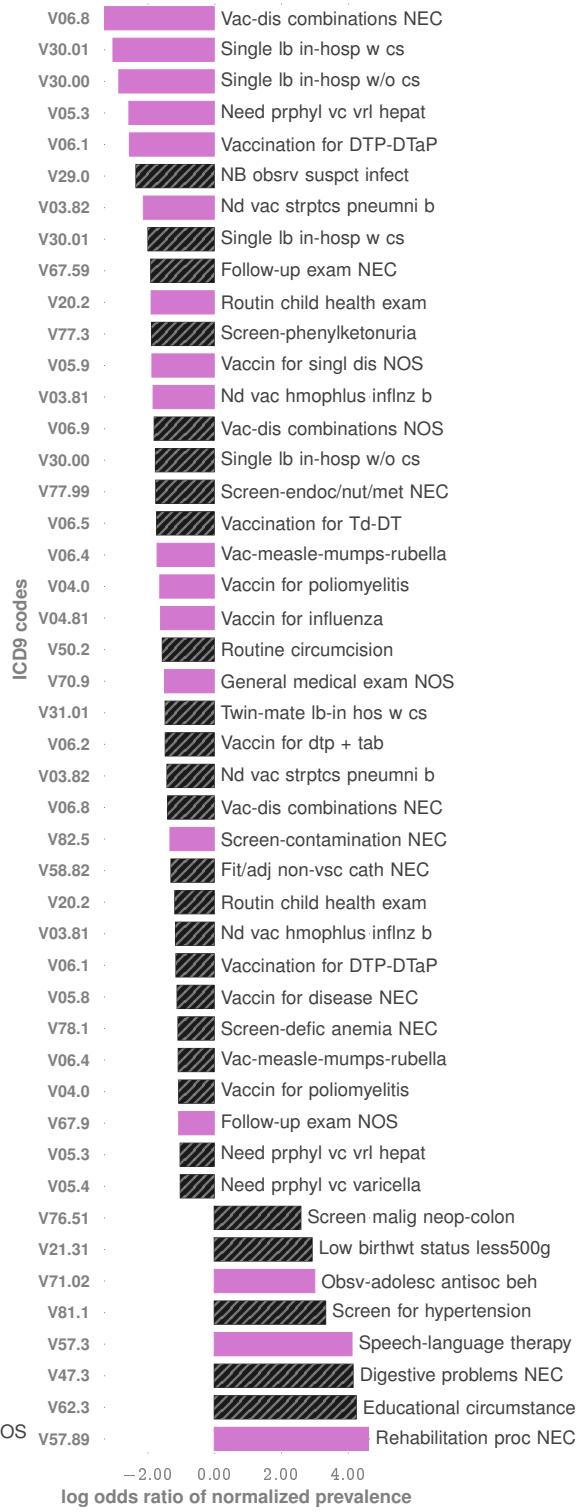
log odds ratio of normalized prevalence

SI-Fig. 1: Details of Co-morbidity Patterns (at age < 3 years) for immunologic (panel A), respiratory (panel B), infections (panel C), and disorders with similar pathobiology manifesting opposing association with autism (panel D).

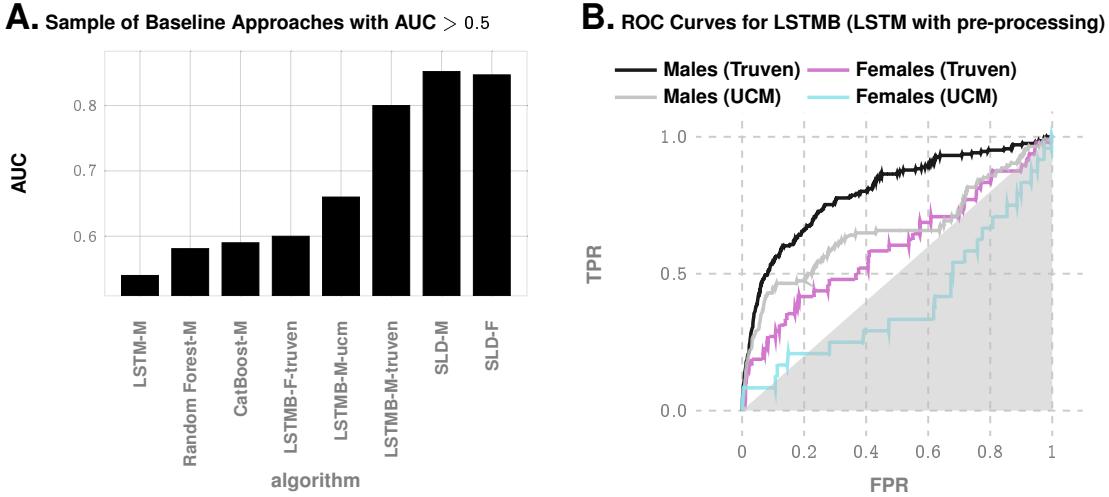
A. Mental Disorders



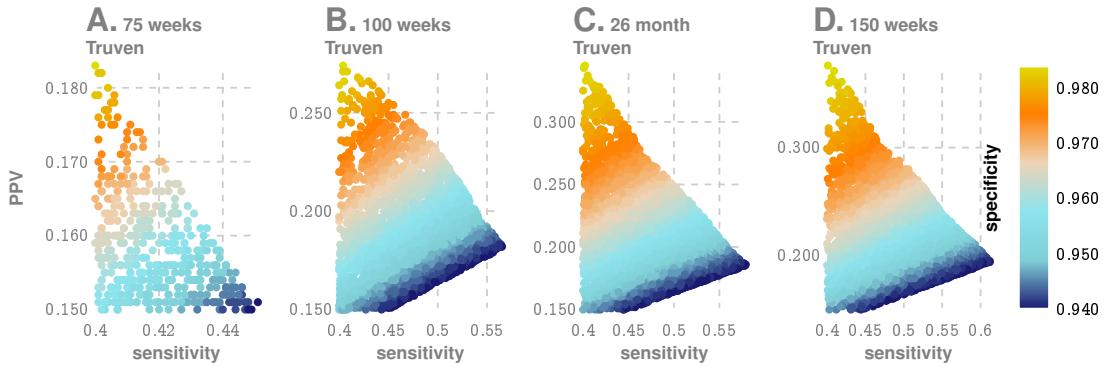
B. Vaccinations & Health Service Encounters



SI-Fig. 2: **Co-morbidity Patterns** for mental disorders, vaccinations and health-service encounters.



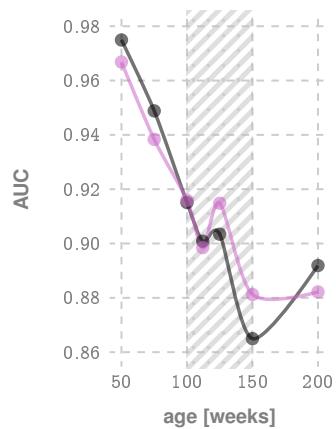
SI-Fig. 3: Performance of standard tools on correctly predicting eventual ASD diagnosis, computed at age 150 weeks of age. Long-short Term Memory (LSTM) networks are the state of the art variation of recurrent neural nets, and Random Forests and Gradient Boosting classifiers (CatBoost) are generally regarded as a representative state of the art classification algorithms. Sequence Likelihood Defect (SLD) is the approach developed in this study. LSTMB denotes LSTM with identical pre-processing as in our pipeline (instead of using raw diagnostic codes). We get much better performance with LSTMB with males in the Truven dataset, but the performance is sensitive to the sizes of the training set, and degrades for smaller samples available for females and in the UCM database, as shown in Panel B.



SI-Fig. 4: **4D Search To Take Advantage of Data on Population Stratification (Using Prevalence of 2.23% as reported by CHOP¹).** While as a standalone tool our approach is comparable to M-CHAT/F at around the 26 month mark (and later), we can take advantage of the independence of the tests to devise a conditional choice of the operating parameters for the new approach. In particular, taking advantage of published estimated prevalence rates of different categories of M-CHAT/F scores, and true positives in each sub-population upon stratification, we can choose a different set of specificity and sensitivity in each sub-population to yield significantly improved overall performance across databases, and much earlier. Additionally, we can choose to operate at a high recall point, where we maximize overall sensitivity, or a high precision point, where we maximize the positive predictive value.

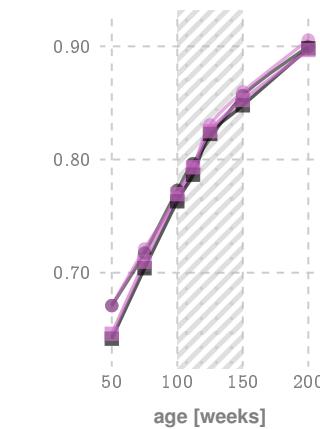
A. Disambiguation of Autism Diagnosis from Other Psych. Phenotypes

—●— Males (Truven)
—●— Females (Truven)



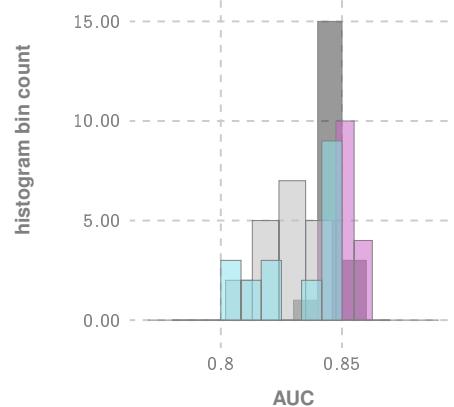
B. Comparison of Performance with One vs Two ASD Diagnostic Codes

—●— Males (Truven, two codes)
—●— Females (Truven, two codes)
—■— Males (Truven, one code)
—■— Females (Truven, one code)

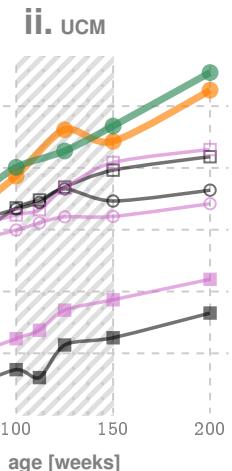
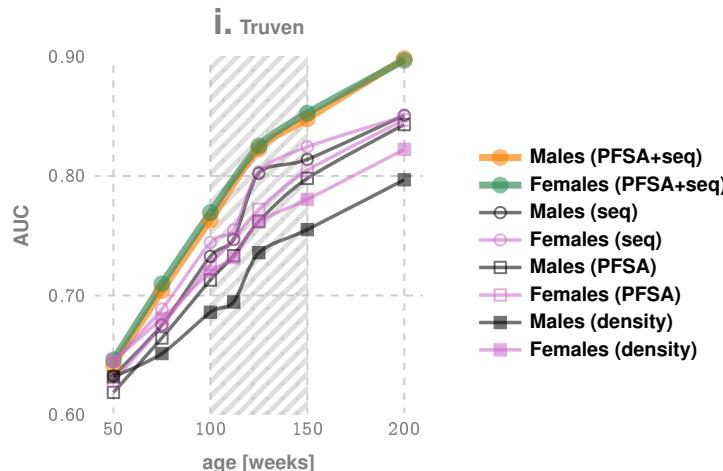


C. AUC Distribution with Matched Control & Treatment Population Sizes

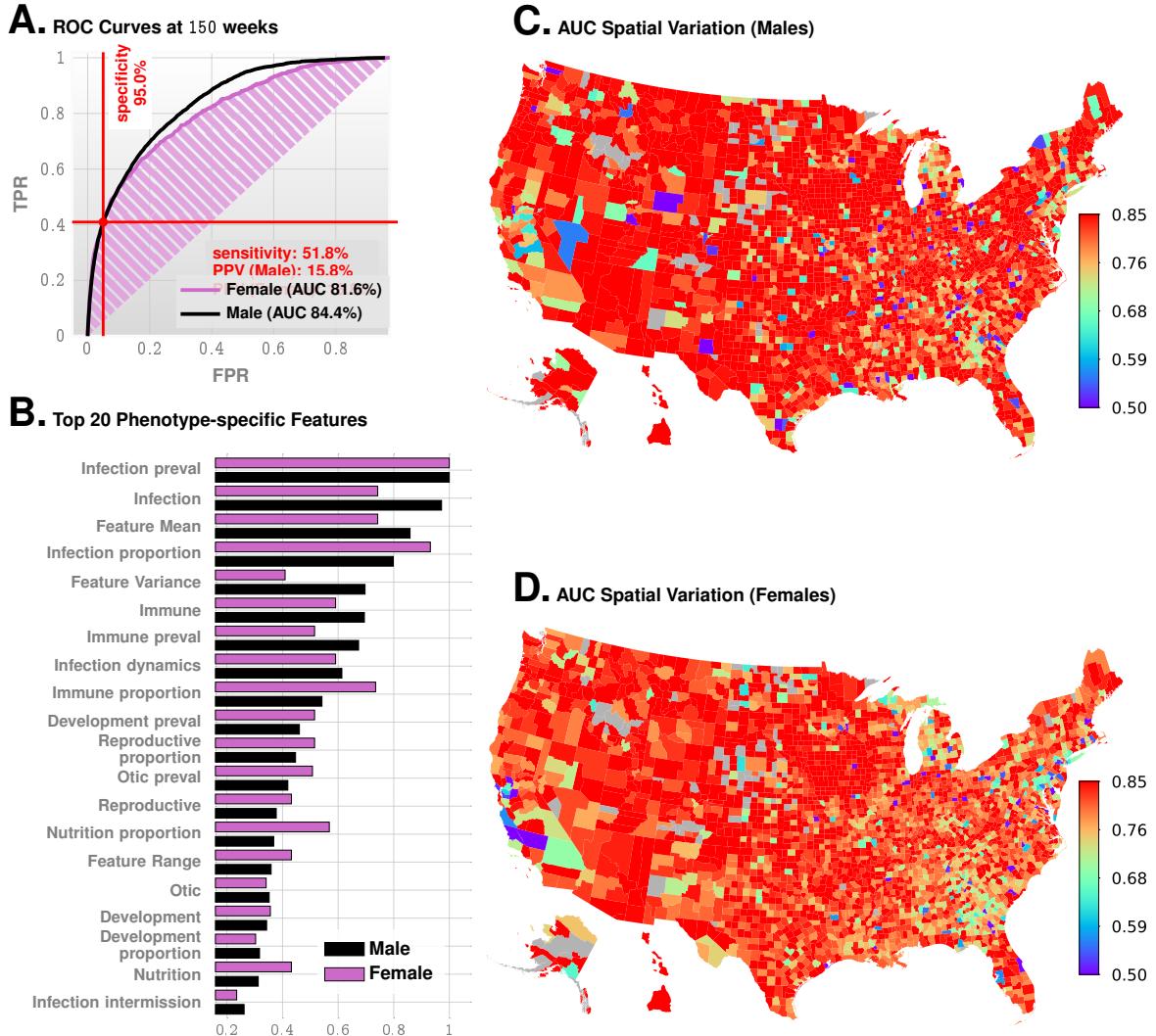
—■— Males (Truven)
—■— Females (Truven)
—■— Males (UCM)
—■— Females (UCM)



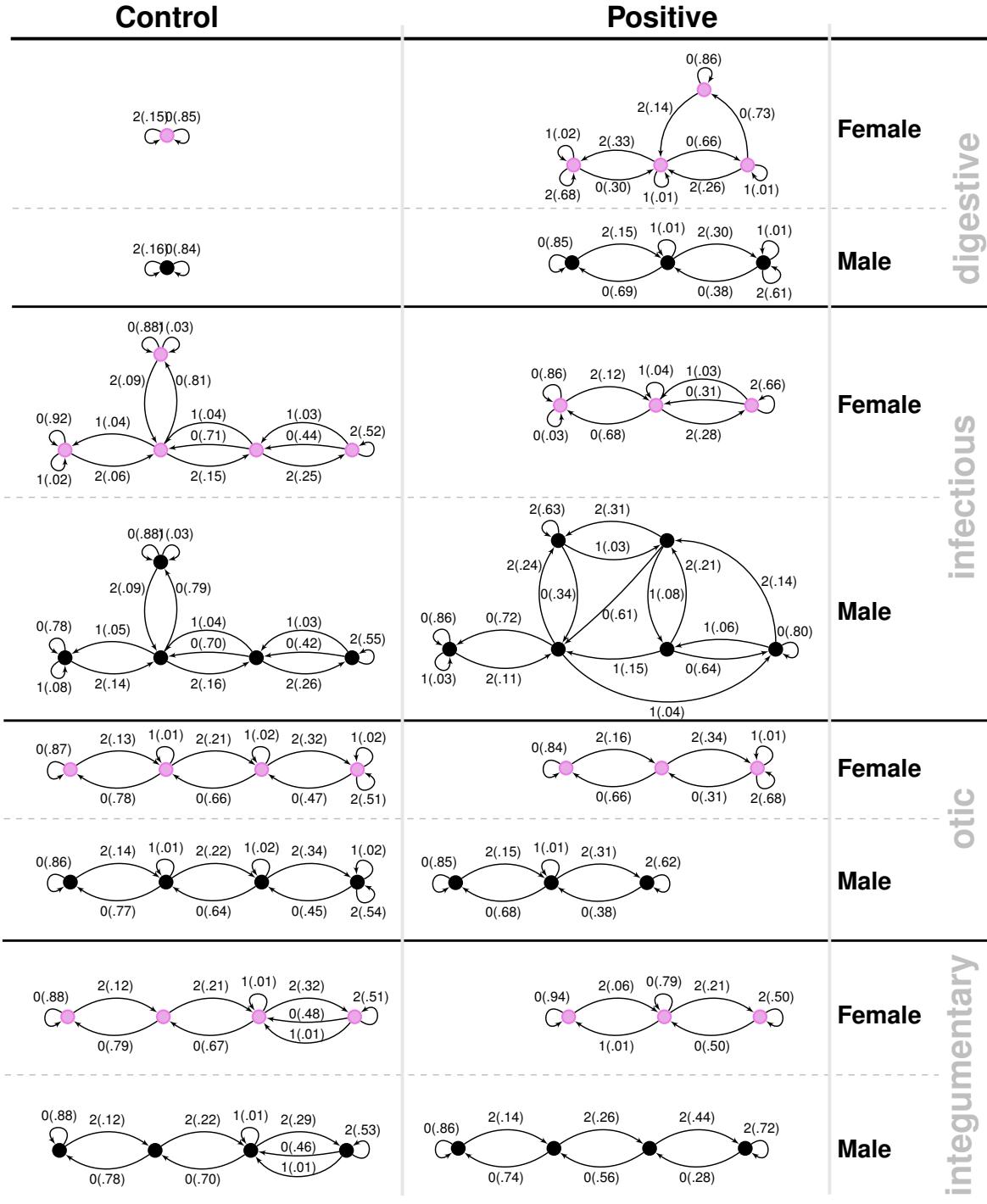
D. Comparison of Performance with Different Feature Categories (Only PFSA based features, Only Sequence-statistics based features, only Code-density, and PFSA + Sequence-statistics features combined)



SI-Fig. 5: Evaluations of Feature Subsets, Class Imbalance, Code Density, Coding Uncertainty, & Disambiguation from Other Psychiatric Phenotypes. Panel A illustrates that the pipeline performance where the control group is restricted to children to have at least one psychiatric phenotype other than ASD. It is clear that we have very good discrimination between ASD and non-ASD phenotypes. Panel B illustrates the situation where we restrict the treatment cohort to children to have at least 2 AD diagnostic codes, to see whether the pipeline performance is markedly different in populations where the coding errors/uncertainty is smaller. We see that such restrictions have no appreciable effect on pipeline performance. Panel C illustrates the AUC distributions obtained by using sampled control cohorts that are of the same size as the treatment cohort, to evaluate the effect of class imbalance. Again we see that such restrictions do not appreciably change performance. Panel D explores the performance changes when we use a restricted set of features, or simply use code density as the sole feature. We conclude that the combined feature set used in our optimized pipeline is superior to using the subsets individually. Code density is the least performant feature, and is not stable across databases.



SI-Fig. 6: Predictive Performance without psychiatric codes (ICD9 290 - 319) and codes for health status and services (ICD9 V0-V91) included. As shown, the performance is comparable at 150 weeks, with the AUC for females marginally lower (compare with Fig. 2 in the main text). The feature importances also are similar, with infectious diseases inferred to have the most importance (or weight) in the pipeline, which is also the case once we add psychiatric phenotypes, and codes for health services in our analysis. As shown in Fig. 2A, the psychiatric codes all increase risk, and the vaccination codes (See Fig. 2B) all decrease risk when those codes are included. This is why an alternate analysis was carried out to make sure that we are not picking up on psychiatric codes alone. Note in particular that the sensitivity/specificity point highlighted in panel A above is identical after adding the codes. This suggests that our predictive performance arises from patterns learned from co-morbidities, which are not just neuropsychiatric in nature.



SI-Fig. 7: Probabilistic Finite State Automata models generated for different disease categories for the control and positive cohorts. We note that in the first cases (digestive disorder), the models get more complex in the positive cohort, suggesting that the disorders become less random. However, in the categories of otic and integumentary disorders, the models become less complex suggesting increased independence from past events of similar nature. In case of infectious diseases, the model gets more complex for males, and less complex for females, suggesting distinct sex-specific responses associated with high ASD risk.

SI-Table 4: Disease Categories With Detailed Set of ICD9 Codes Used

Cat. †	Description	Constituent ICD9 Codes
Hematologic	Diseases Of The Blood And Blood-Forming Organs	286.9 286.7 286.6 283.19 283.10 283.11 283.9 283 283.1 284.9 284.8 284.81 284.0 284.89 284.09 284 284.01 282.2 287.49 287.41 287.39 287.4 287.5 287.32 287.3 287.30 287.31 286.3 286.2 286.1 286.0 286.4 282.1 282.6 282.5 282.41 282.42 282.68 282.69 282.62 282.63 282.60 282.61 282.64 282 282.8 287.33 281.2 281.3 280.0 282.9 285.8 285.9 280.9 284.2 285.1 285.2 285.3 280.1 285.22 285.21 282.3 776.5 285 283.0 285.29 280.8 282.7 282.40 282.49 284.1 284.19 284.12 284.11 281.8 281.9 281.4 281.0 281.1 286.5 287 287.8 287.9 287.2 287.0 287.1 285 289.52 289.50 289.51 289.59 289.4 289.5 289.81 289.83 289.82 289.89 289 289.7 289.8 289.9
Psychiatric	Mental Disorders (Except ASD)	290 through 319 (except 299.x)
Metabolic	Metabolic Disorders (Distinct from respiratory, digestive and immunological conditions)	273.4 270 270.2 270.3 712.11 712.10 712.13 712.12 712.15 712.14 712.17 712.16 712.19 712.18 712.31 712.30 712.37 712.36 712.35 712.34 712.38 712.33 712.32 712.28 712.29 712.24 712.25 712.26 712.27 712.20 712.21 712.22 712.23 712.39 712.1 712.3 712.2 277.6 275.1 277.5 277.87 270.7 270.6 276.6 276.4 276.2 276.3 276.0 275.41 276.1 276.8 276.9 276.6 275.5 275.42 271.1 330.2 272.7 271.2 274.81 274.85 274.81 274.82 712.99 712.98 274.01 274.04 274.03 274.02 712.91 712.90 712.93 712.92 712.95 712.94 712.97 712.96 712.88 712.89 274.10 274.11 712.82 712.83 712.80 712.81 712.86 712.87 712.84 274.19 712.9 712.8 274.0 274.1 274.2 274.8 274.9 271.2 275.01 270.5 270.4 278.8 272.3 275.03 275.09 271.3 272.6 272.5 278.1 271.8 277.5 263.0 263.2 262 260 261 263 263.1 269.8 269.9 263.8 263.9 269 277.7 272.2 330.3 271.9 275.40 272.8 277.8 275.49 275.2 277.8 275.4 269.3 275.9 275.8 277.9 277.89 251.2 251.1 251.0 278.01 278.01 278.03 270.8 270.9 278 278.0 278.02 277.86 270.1 275.3 277.1 277.81 277.82 272 272.1 277.2 272.4 272.9 273.9 273.8 268.1 265.2 268.0 268.2 268 265.0 265.1 266.1 266.0 266.2 266.9 264.3 264.2 264.1 264.0 264.7 264.6 264.5 264.4 264.9 264.8 268.9 267 266 265 264 269.2 269.0 269.1 278.2 278.3 278.4
Cardiovascular	Diseases Of Arteries, Arterioles, And Capillaries	442.89 441.6 442.84 442.82 442.83 442.8 441.03 441.02 441.01 441.00 441 414.19 414.12 442 414.10 414.11 447.70 447.71 447.72 447.73 414.1 442.81 441.9 442.1 442.0 442.3 442.2 441.2 441.3 441.0 441.1 442.9 441.7 441.4 441.5 437.3 447.7 443.29 443.23 443.22 443.21 443.24 443.2 444.2 444.8 444.81 444.82 444.1 444.0 444.0 444.89 444.89 444.22 444.21 445.81 440.31 440.30 440.32 444.01 414.00 414.03 414.02 414.05 414.04 414.07 414.06 445.89 411.81 445.02 445.01 440.24 440.22 440.23 440.20 440.21 440 445.40 429.40 414.0 414.2 440.4 440.3 440.2 440.1 440.0 440.9 440.8 445.8 445.0 414.3 414.4 426.54 426.53 426.52 426.51 426.50 426.13 426.12 426.11 426.10 426.89 426.9 426.8 426.81 426.3 426.2 426.1 426.0 426.7 426.6 426.5 426.4 427.61 427.60 427.5 427.89 427.69 427.32 427.41 427.9 427.81 427.8 427.42 427.6 427.4 427.3 427.1 427.2 427.3 427.0 427.1 425.8 425.9 425.4 425.7 425.0 425.1 425.2 425.3 346.6 438.51 438.52 438.53 438.50 290.4 431.0 431.0 438.42 438.41 438.40 432.9 432.9 290.42 290.41 290.40 432.0 432.1 433.00 433.01 434.9 346.61 346.60 346.63 346.62 433.80 433.81 433.11 433.10 438.32 438.30 438.31 434.0 433.91 433.90 430.0 430.10 434.11 434.1 438.21 438.20 438.22 433.20 433.21 438.6 438.7 438.4 438.5 438.2 438.0 438.1 438.8 438.9 438.43 438.42 438.40 434.91 434.01 434.00 438.10 438.11 438.12 438.13 438.14 438.19 433.31 433.30 438.85 438.84 438.83 438.82 438.81 438.89 433.9 433.8 433.1 433.0 433.2 433.2 437 435.3 435.1 435.0 435.9 435.8 437.5 437.4 437.7 437.6 437.1 437.0 431 437.9 437.8 459.3 459.33 459.32 459.31 459.30 459.39 452 453 453.7 453.6 453.5 453.4 453.3 453.2 453.1 453.0 453.9 453.8 453.52 453.51 453.50 453.79 453.71 453.73 453.72 453.75 453.74 453.77 453.76 453.84 453.89 453.40 453.41 453.42 453.81 453.82 453.83 415.11 453.85 453.86 453.87 405.0 405.1 404.9 403.11 402.00 402.01 404.1 404.0 402.1 402.0 403.0 405.99 405.91 402.9 402.91 402.90 405.11 401.0 401.1 404.00 404.01 404.02 404.03 405.19 401.9 405 404 403 402 401 405.9 403.01 403.00 402.11 402.10 403.10 404.13 404.12 404.11 404.10 405.01 403.9 405.09 403.90 437.2 403.1 403.91 404.93 404.92 404.91 404.9 448 448.5 458.0 458.2 458.1 458.9 458.8 458.9 458.21 426.82 429.71 410.01 410.00 410.02 410.41 410.40 410.40 410.42 410.22 410.21 410.20 429.7 410.70 410.71 410.72 429.79 410.92 410.90 410.91 410.30 410.31 410.32 410.12 410.10 410.11 410.52 410.50 410.51 410.41 410.6 410.7 410.0 410.1 410.2 410.3 410.8 410.9 411.0 410.62 410.61 410.60 410.41 410.42 410.81 410.40 410.82 424.1 424.0 424.2 424.3 424.2 429.89 429 429.1 429.5 429.6 429.8 429.9 459 276.5 429.2 429.3 428.9 428.4 428.1 428.0 428.3 428.2 429.81 429.83 429.82 428.32 428.31 428.30 428 459.8 459.9 459.0 276.50 276.51 428.42 428.43 428.40 428.41 276.52 428.20 428.21 428.22 428.23 428.00 459.89 448.1 454 455 455.9 455.8 454.8 454.9 455.1 455.0 455.3 455.2 455.5 455.4 455.7 455.6 454.2 447 454.0 454.1 454.7 357.32 447.8 447.9 448.9 447.4 447.5 447.0 447.2 447.3 414.1 413.1 413.0 413.9 411.89 411.1 411 414.9 413.8 414.8 411.8 443.89 443.8 443.9 443.81 459.81 443.0 416.2 415.19 415.1 415.13 415.12 416.1
Endocrine	Disorders Of Thyroid and other Endocrine Glands	244 244.9 244.8 244.2 255.41 255.42 255.5 255.4 255.2 255.1 255.13 259.2 243 255 253.5 259.4 255.11 242.2 240.0 241 240 240.9 241.0 242.20 242.21 241.9 241.1 253.3 704.1 255.12 255.10 255.14 246.9 246.8 246.3 246.1 246.0 246 255.3 255.9 255.8 255.6 255 252.8 252.9 252.0 252.1 252.01 252.02 252 252.08 259.52 253.4 253.6 253.1 253.0 253.2 253 253.9 253.8 242.1 242.10 242.4 242.9 242.8 242.11 242.40 242.41 376.32 242 242.31 242.30 242.91 242.90 242.80 242.81 250.20 250.30 250.22 250.42 250.40 250.02 250.00 250.12 250.82 250.80 250.90 250.92 250.50 250.52 250.72 250.60 250.70 250.62 362.05 362.01 362.07 362.06 362.03 357.2 362.03 250.7 362.02 250.8 250.9 362.04 250.2 250.3 250 250.1 250.6 250.0 250.4 250.5 259 259.3 259.50 259.9 259.8 258.8 258.9

† Categories inferred to be important for risk modulation are highlighted. Continued on next page

Integumentary	Diseases Of Skin And Subcutaneous Tissue	706.0 706.1 704.0 704.00 704.02 704.01 704.09 680.9 680.8 680.1 680.0 680.3 680.2 680.5 680.4 680.7 680.6 680 698.8 698 698.1 698.0 698.9 757.31 757.5 757.4 757.33 700 694.2 694.3 694.0 694.1 709.0 00 709.01 709.09 757.39 757.3 695.59 695 695.57 695.56 695.55 695.54 695.53 695.52 695.51 695.50 695.5 695.58 704 704.9 704.2 704.3 704.8 757 757.1 702.11 702.19 702.0 702.1 697.1 697.0 697.9 697 697.8 759.82 703.8 703.9 703.0 703.80 703 757.8 757.9 757.2 525.9 525.8 525.5 525.4 525.7 525.6 525.1 525.0 525.3 525.79 525 525.73 525.72 525.71 521.35 521.24 521.25 521.20 521.21 521.22 521.23 521.42 521.40 521.41 521.49 522.2 522.3 522.0 522.1 522.8 522.9 521.33 521.32 521.31 521.30 525.19 521.522 521.34 525.11 525.10 525.13 525.12 521.89 525.42 525.43 525.40 525.41 525.44 521.08 521.09 521.02 521.03 521.00 521.01 521.06 521.07 521.04 521.05 521.1 521.0 521.3 521.2 521.5 521.4 521.7 521.6 521.9 521.8 525.51 525.50 525.52 525.54 521.11 521.10 521.13 521.12 521.15 521.14 521.81 525.69 525.64 525.65 525.66 525.67 525.60 525.61 525.62 525.63 523.11 523.10 523.31 523.30 523.33 523.32 523.3 522.7 522.4 522.5 523.6 523.5 523.4 523.9 523.8 523 522.6 523.00 523.01 523.2 523.1 523.40 523.41 523.42 523.0 523.22 523.23 523.20 523.21 523.24 523.25 709.1 707.00 707.02 707.03 707.04 707.05 707.06 707.01 709.9 707.09 707 701 702 707.25 707.2 707.1 707.0 707.22 707.8 707.9 707.23 706.8 707.20 702.8 705.2 706.2 707.07 701.9 701.8 709.3 709.8 701.1 701.0 701.3 701.2 701.4 707.24 707.19 707.15 707.14 707.13 707.12 707.11 710 705.21 705.22 707.21 707.10 710.9 710.8 709 695.9 709.2 729.7 729.72 729.91 729.92 729.9 729.79 729.99 729.0 729.71 729.73 729.72 528.8 528.9 528.4 528.5 528.6 528.7 528.0 528.1 528.2 528.0 528.79 527.3 528.09 528.3 528 528.71 528.72 692.7 692.71 692.72 692.72 692.82 692.74 692.75 692.76 692.77 692.73 692.79 558.1 508.0 508.1 655.6 595.82 655.63 655.60 655.61 737.33 737.11 695.3 706 706.9 690 690.1 690.10 690.11 690.12 690.8 690.18 705.9 705.1 705.0 705.89 705 705.82 705.8 529.9 529.8 529.1 529.0 529.3 529.2 529.5 529.4 529.6 529
Respiratory	Diseases Of The Respiratory System (Distinct from Infectious)	516.31 516.30 516.33 516.32 516.35 516.34 516.37 516.36 516.9 516.8 516.5 516.3 516.2 516.1 516.0 516 277.00 277.01 277.02 277.03 277.09 470 478.0 491.9 491.8 515 491.1 491.0 491.2 515.0 491.22 491.21 494.0 492.0 492.8 494.1 494 496 491.20 491 492 478.32 478.33 478.30 478.31 478.34 478.75 478.4 478.79 478.3 478.74 478.5 478.6 478.7 478.70 478.71 471.8 471 471.9 471.1 471.0 519.09 519.00 519.02 519.02 495.2 495.5 495.4 495.7 495.6 495.1 495.0 495.3 518.4 495.4 495.9 495.8 517 503 519.8 507.0 517.8 504 505 502 517.3 500 501 518.84 518.82 518.83 518.81 514.0 519.11 518.89 519.0 519.4 518.8 518.5 519.9 518.1 518.0 518.3 518.2 501.0 512.1 512.0 517.2 512.8 519 518 510.0 510 512 510.9 514 500.0 327.24 327.25 327.26 327.2 327.20 327.21 327.22 327.23 327.27 327.29
Digestive	Diseases Of The Digestive System	540.0 540.1 540.9 541.0 542 543 540 541 543.9 543.0 562 562.1 562.11 562.03 562.02 562.01 562.00 562.10 562.0 562.12 562.13 530.3 530.13 530.11 530.0 530.6 530.5 530.4 456.21 530.9 530.8 456.2 456.20 456.0 456.1 530.89 530 530.84 530.85 530.81 530.82 530.83 560.89 579 560.30 560.0 564.09 560.1 564.01 564.00 564.02 564.81 560.39 564.89 560.31 560.32 579.9 536.8 564 537.0 537.1 537.2 537.3 536.9 537.5 537.6 564.5 537.8 537.9 560.2 560.3 536.1 536.0 536.3 537.81 560.8 560.9 537.89 537 536 564.8 564.9 537.4 564.7 564.0 538 532.51 532.50 532.53 532.21 533.20 569.41 534.11 534.10 531.70 531.71 534.1 531.9 534.9 532.6 532.61 533.50 533.51 534.60 534.61 531.61 531.60 531.6 531.7 531.4 531.5 531.2 531.3 531.0 531.1 531.1 531.9 534.30 532.53 532.51 531.51 533.41 533.40 532.71 532.70 534.71 534.70 530.20 530.21 534.9 531.41 531.40 532.00 532.01 534.40 534.41 532.11 532.10 531.30 531.31 534.51 534.50 569.89 533.70 533.71 569.82 533.10 533.11 532.1 532.0 532.3 532.2 531.21 531.20 532.7 532.6 532.9 532.91 532.90 533.30 533.31 532.40 532.41 534.7 534.6 534.5 534.4 534.3 534.2 534.1 534.0 533.4 533.5 533.6 533.7 533.0 533.1 533.2 533.3 533.9 531.01 531.00 530.2 534.00 530.41 532.31 532.30 532.5 532.4 751.3 568.0 578.9 569.9 568.9 558.8 568.8 568.9 568.81 568.82 568.89 251.5 251.4 577 251.9 251.8 577.8 577.9 577.8 577.9 579.4 251 577.0 577.1 527.2 569.4 569.1 565 566 569.2 569.42 569.43 566.0 569.44 565.1 565.0 564.6 569.49 527.9 527.8 527.7 527.6 527.5 527.4

[†] Categories inferred to be important for risk modulation are highlighted. Continued on next page

Developmental	Congenital anomalies (Non-overlap. with musculoskeletal)	<p>755.55 743.45 743.12 743.11 743.10 743.06 743.00 743.03 743.1 743.44 743.41 743.42 743.43 743.22 743.2 743.20 743.21 758.4 745.1 745.0 745.3 745.2 745.5 747.5 746.09 747.745.4 745.8 746.02 746.01 746.00 745.9 745.7 747.89 745.7 745.6 747.82 747.83 745.60 745.61 745.69 747.21 746.6 746.7 746.4 746.5 746.2 746.3 746.0 746.1 747.9 747.8 746.8 746.9 746.87 746.86 746.85 746.84 746.83 746.82 746.81 746.89 745.19 747.11 745.11 745.10 745.12 756.4 744.4 743.63 743.61 743.66 744.1 743.64 743.65 743.69 744.8 744.9 744.0 754.1 754.0 748.1 748.0 744.744.49 744.2 744.84 744.3 744.5 743.6 744.29 744.47 744.42 744.89 744.09 744.81 744.83 744.82 744.05 744.04 744.01 744.00 744.03 744.02 744.41 744.43 744.24 744.23 744.22 744.21 744.46 743.35 742.51 756.2 756.19 742.53 756.11 756.10 756.13 756.12 756.15 756.14 758.31 315.4 315.5 315.8 315.9 758.2 752.4 752.1 752.0 752.3 752.2 752.42 752.43 752.40 752.41 752.47 752.44 752.45 752.49 752.36 752.35 752.34 752.33 752.32 752.31 752.39 752.19 752.11 752.10 759.83 751.61 751.7 751.6 751.60 751.69 743.62 520.4 520.5 520.0 520.1 520.2 752.7 758.7 756.16 752.5 752.6 752.51 752.52 752.69 752.64 752.65 752.61 752.62 752.63 740 740.0 740.2 741.9 740.1 741.02 741.0 741.00 741.01 741.742.0 742.0 742.1 742.4 742.5 742.9 524.70 524.71 524.72 524.73 524.74 524.75 524.76 524.79 520.6 520.7 520.3 520.8 520.9 750.29 750.26 750.27 750.25 750.22 750.23 750.21 524.2 524.3 524.4 524.5 524.7 524.8 524.9 524.5 520 749.04 749.02 749.03 749.00 749.01 524.39 750.16 750.15 750.13 750.12 750.11 750.10 750.19 749.14 749.11 749.10 749.13 749.12 524.81 524.82 524.89 749.20 749.21 749.22 749.23 749.24 749.25 524.34 524.35 524.36 524.37 524.30 524.31 524.32 524.33 750.2 750.1 750.0 524.56 524.57 524.54 524.55 524.52 524.53 524.50 749.524.59 749.1 749.0 749.2 524.23 524.22 524.21 524.20 524.27 524.26 524.25 524.24 524.29 524.28 758.1 759.81 259.1 259.0 748.61 748.60 748.748.69 748.8 748.9 748.2 748.3 748.4 748.5 748.6 758.33 759.89 758.6 313.89 313.9 313.309.21 307.7 307.6 313.8 313.23 313.2 759.9 759.2 759.3 759.0 759.1 759.7 759.4 747.40 747.41 747.64 747.29 747.62 747.63 747.60 747.61 747.22 747.49 747.20 747.81 747.3 747.2 747.1 747.0 747.6 747.4 747.42 747.10 758.32</p>
Nutrition	Nutrition, metabolism, and development	7830,78321,7833,78340,78342,7837,7839
Health Status & Services Contact.	Vaccination, Supplementary Classification Of Factors Influencing Health Status And Contact With Health Services etc.	V01-V91
		† Categories inferred to be important for risk modulation are highlighted.

Algorithm 1: ICD-9 Encoding

```

input : Dataset, TargetDiseaseGroup, DiseaseGroups
output: Encoding

1 Encoding ← new Dictionary();
2 for diseaseGroup ∈ DiseaseGroups do
3   Encoding[diseaseGroup][patientID] ← new List();
4   Encoding[diseaseGroup][gender] ← new List();
5   Encoding[diseaseGroup][record] ← new List();
6   Encoding[diseaseGroup][target] ← new List();
7   for record ∈ Dataset do
8     //encode Dataset into a weekly trinary sequence;
9     weeklyEncoding ← new List();
10    for weeklyDiseaseRecord ∈ record do
11      //no code recorded for the observed week;
12      if weeklyDiseaseRecord.code == NIL then
13        | append "0" to weeklyEncoding;
14      if weeklyDiseaseRecord.code ∈ diseaseGroup.codes then
15        | append "1" to weeklyEncoding;
16      if weeklyDiseaseRecord.code ≠ diseaseGroup.codes then
17        | append "2" to weeklyEncoding;
18    target ← 1 if any weeklyDiseaseRecord.code of record ∈ TargetDiseaseGroup;
19    if target == 1 then
20      | cut weeklyEncoding up to (but not including) first occurrence of TargetDiseaseGroup member;
21      append record.patientID to Encoding[diseaseGroup][patientID];
22      append record.gender to Encoding[diseaseGroup][gender];
23      append weeklyEncoding to Encoding[diseaseGroup][record];
24      append target to Encoding[diseaseGroup][target];
25 return Encoding;

```

1. Pipeline Optimization

1.1. Input Data Format

To encode the ICD-9 codes, 17 Disease Groups of codes are used to transform the raw health records into a format suitable for PFSA. As described in *Algorithm 1*, for each patient, the list of ICD-9 codes is encoded into a weekly array of three-symbol alphabet digits with respect to selected disease group, for each week: "0" - no disease "1" - disease from the selected group, "2" - other disease.

Once the trinary encodings are ready, the PFSA pairs are fit for each of the disease groups, on positive (treatment) and negative (control) sets using genESeSS algorithm² (See Section 11), as described in *Algorithm 2*. The PFSA pairs are then used to obtain the loglikelihood scores of belonging to a PFSA modeling the positive and the control cohorts accordingly for each of the encodings of a patient record. As a result, we yield the difference between positive and control loglikelihoods for each disease group of each patient. The positive value of difference means that with respect to a given disease group, a certain patient is more likely to be a positive one. Conversely, the negative value of difference signifies that a patient is more likely to be from the control group. These features, as well as their aggregations and the aggregations of the ternary encoding arrays, are used as the features for the final LightGBM gradient boosting classifier.

1.2. Algorithms

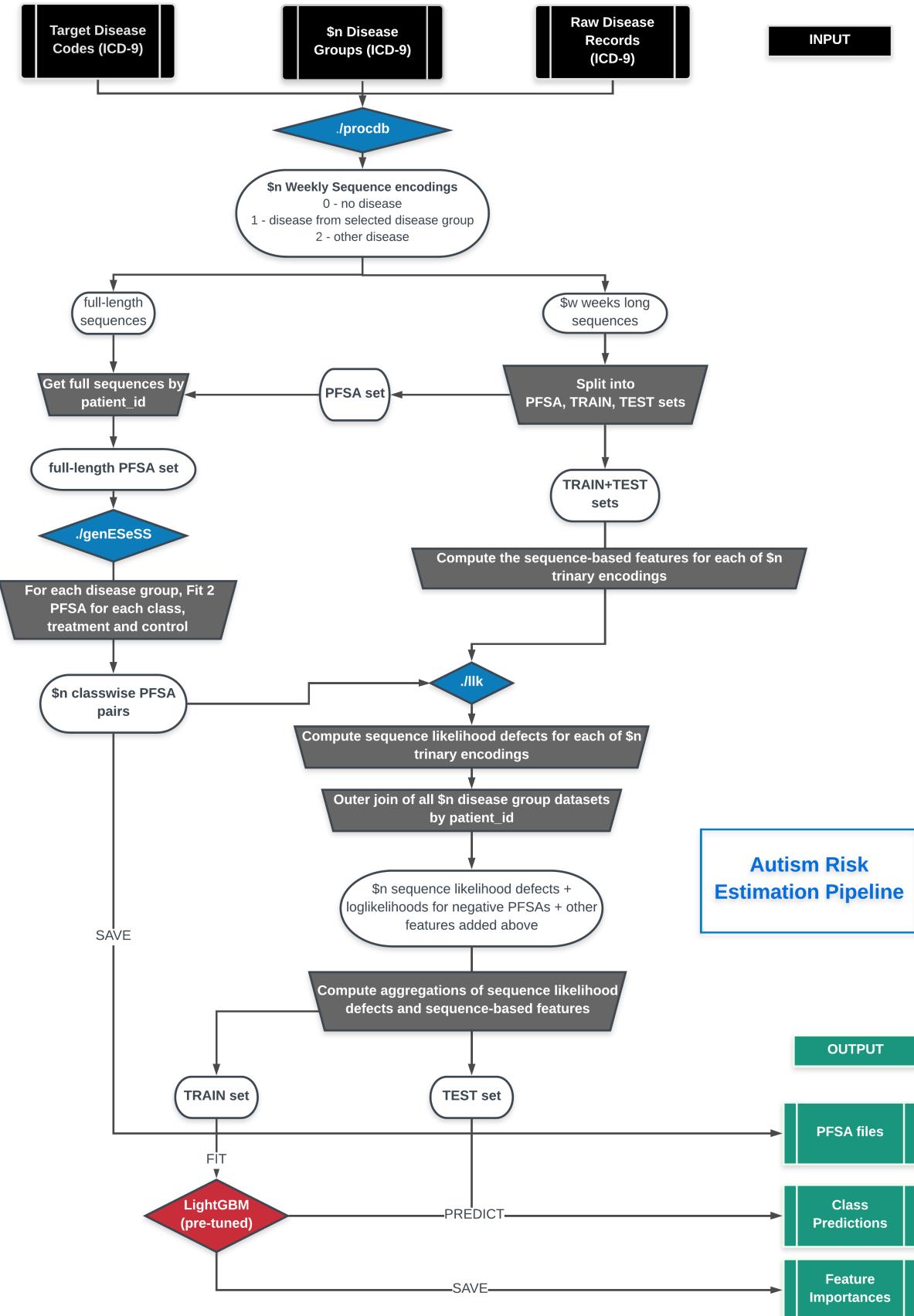
The key data processing approach is outlined in Algorithm 1. The remaining steps of the approach are sketched in Algorithm 2. Fig. 8 shows the overall schema, including the breakdown of a database into a test set, and two training sets: one for training the HMM models, and one for training the boosting classifier.

Algorithm 2: Prediction Pipeline Training

input : Encoding, DiseaseGroups, SequenceFeatures, hyperparameters
output: Predictions, FeatureImportances

```

1 DiseaseDatasets ← new Dictionary () ;
2 for Dataset, DiseaseGroup ∈ zip(Encoding, DiseaseGroups) do
3   PFSAsset, LLset ← TrainTestSplit (Dataset, w.r.t = "target") ;
4   df ← new Dataframe () ;
5   df[patientID] ← LLset[patientID];
6   df[target] ← LLset[target];
7   //Generate 2 PFSA for each class;
8   PositivePFSAsset ← PFSAsset[PFSAsset.target == 1];
9   NegativePFSAsset ← PFSAsset[PFSAsset.target == 0];
10  PosPFSA ← genESeSS (PositivePFSAsset) ;
11  NegPFSA ← genESeSS (NegativePFSAsset) ;
12  //For each record, compute loglikelihoods of being generated by either of 2 PFSA generated above;
13  PosLLK ← llk (LLset, PosPFSA) ;
14  NegLLK ← llk (LLset, NegPFSA) ;
15  //Compute sequence likelihood defect;
16  df[DiseaseGroup] ← pairwise (PosLLK - NegLLK) ;
17  df[DiseaseGroup + '_abs_neg'] ← NegLLK;
18  for SequenceFeature ∈ SequenceFeatures do
19    df[DiseaseGroup + '_' + SequenceFeature] ← [ComputeSequenceFeature (SequenceFeature,
      seq) for each seq ∈ LLset['record']];
20  DiseaseDatasets[DiseaseGroup] ← df;
21 Dataset ← outerjoin (DiseaseDatasets.values, on = 'patientID') ;
22 Aggregate all features in Dataset where feature_name ∈ DiseaseGroups (mean, std. deviation, range);
23 Aggregate all features in Dataset where feature name minus '_abs_neg' ∈ DiseaseGroups (mean, std.
  deviation, range);
24 Aggregate all sequence features in Dataset (mean, std. deviation, range, max);
25 TrainSet, TestSet ← TrainTestSplit (Dataset, w.r.t = "target") ;
26 LGBM ← new LightGBM(hyperparameters) ;
27 LGBM.fit(TrainSet);
28 Predictions ← LGBM.predict(TestSet);
29 return Predictions, LGBM.feature_importances;
```



SI-Fig. 8: Pipeline schema: How the data set is split into test sets and two training sets: one for inferring HMM models, and one for training the boosting classifier. The two ket algorithms here are `genESeSS`² and the `llk` which does the sequence likelihood computation described in Section 12

2. Example Run with Released Application

Navigation

Project description

downloads 174/month

pypi v1.0.30

Last released: May 8, 2019

Zero-Knowledge Risk Oracle for predictive diagnoses of childhood neuropsychiatric disorders from sparse electronic health records

Project description

Release history

SI-Fig. 9: Screen capture of the page on pypi.org hosting the released application Link: <http://pypi.org/ehrzero>

```
[1]: from ehrzero import ehrzero as ehr
      import warnings
      warnings.filterwarnings("ignore")

[2]: source = 'test_free.dat'
      outfile = 'out.dat'
      first_weeks = [200, 100] # number of first weeks of the observations to consider
      risks = ehr.predict_with_confidence(source,
                                           outfile,
                                           separator = ',',
                                           delimiter = '|',
                                           n_first_weeks = first_weeks)

[3]: risks
```

	patient_id	week	risk	relative_risk	confidence
0	AAAbby	200	0.000174	0.028200	0.920977
0	AAAbby	100	0.000135	0.021954	0.954741
1	ALorax	200	0.000101	0.016416	0.989583
1	ALorax	100	0.000099	0.016071	0.986710

SI-Fig. 10: Python code prediction example

2.1. Prerequisites & Installation

The minimum prerequisites for running ehrzero are the following:

1. A x64 system running any flavor of Linux.
2. A working python 3.x installation
3. scikit-learn, version = 0.20.0

Installation:

pip3 install ehrzero --user

2.2. EHR data format

Diagnostic data stored in text file, one line per patient as follows: patient id, gender, and list of space-separated, comma-delimited diagnosis records, all separated by spaces. Each diagnosis record consists of the week since the start of the observation, followed by a comma, and the ICD-9 code of the diagnosis.

Example of a patient line:

```
Lorax,M 5,277.03 10,611.79 18,057.8 58,157.8 78,057.8 108,057.8 128,057.8 148,057.8
```

2.3. Sample Python code risk estimation

Once the patient diagnostic data is in the required format, for function `predict_with_confidence` we specify the filepath of the data and the list of the cutoffs for the first weeks since the start of observations for the data we want to analyze. We also specify the separator and delimiter for the patients within file (space and comma are default values, but can be changed for user convenience).

The `predict_with_confidence` function returns the predicted risk of autism for every patient in the input file with all the specified numbers of first weeks to consider.

2.4. Sample Python script risk estimation

The script version is similar to the one mentioned before.

Once `ehrzero` package is installed, locate its directory and go to `texttt../ehrzero/example`. Select one of the `".dx"` or `".dat"` files in `/ehrzero/example/tests` as input and run the following command as an example:

```
python zero.py -data tests/ZEROexample.dat -outfile predictions.csv -nweeks 100 200 300 -Verbose 1
```

```
(base) [onishchenko@midway2-login1 example]$ cat tests/ZERO_example.dat
M:44 380.10:101 381.81:111 084.6
M:9 380.10:104 381.81:11 084.6
M:9 380.10:104 381.81:11 084.6:98 380.11
M:9 380.10:104 381.71:11 084.6
M:9 390.11:4 390.11:4 381.71:11 084.6
(base) [onishchenko@midway2-login1 example]$ python zero.py -data tests/ZERO_example.dat -outfile predictions.csv -n_weeks 100 200 300 -Verbose 1
patient_id week risk relative_risk confidence
A00000001 100 0.001703 0.098002 68.54
A00000001 200 0.001834 0.105517 72.43
A00000001 300 0.001834 0.105517 72.43
A00000002 100 0.002039 0.116816 60.94
A00000002 200 0.001426 0.082033 80.01
A00000002 300 0.001426 0.082033 80.01
A00000003 100 0.001703 0.098002 68.54
A00000003 200 0.001766 0.098170 74.53
A00000003 300 0.001853 0.106625 72.18
A00000004 100 0.000512 0.029449 94.86
A00000004 200 0.001450 0.083436 79.53
A00000004 300 0.001450 0.083436 79.53
A00000005 100 0.000658 0.037868 91.73
A00000005 200 0.000658 0.037868 95.21
A00000005 300 0.000658 0.037868 95.21
(base) [onishchenko@midway2-login1 example]$
```

SI-Fig. 11: Python script prediction example

3. Comparison With State of the Art Off-the-shelf ML Algorithms

Off the shelf algorithms with little or no pre-processing, *i.e.*, using the diagnostic codes themselves are timestamped categorical features failed to produce clinically relevant performance (See Fig. 3). Classifiers such as random forests³, and gradient boosters⁴ might be penalized due to their inability to take into account long-range temporal information. Since the number of diagnostic codes available per patient is small, recurrent neural network implementations such as LSTM⁵ might be suffering from the data sparsity in training. It is possible that the performance of the competing approaches might be improved with extensive tuning or clever feature-engineering.

4. Comparison With Pipeline Variations, Feature Subsets and Neural Net Post-processing

In addition to the naive baseline approaches, we also evaluated the performance achievable with LSTMs (denoted as LSTM_B in Fig. 3) that use identical preprocessing as our pipeline, *i.e.*, representation of diagnostic histories as trinary sequences in 18 categories for each patient, and achieved $\sim 80\%$ AUC at 150 weeks for males in the Truven database (compared to $> 85\%$ for our approach). However, the performances drop significantly when the number of positive samples is reduced, yielding an AUC of 66% on the UCM dataset for males, 60% for females on the Truven dataset, and a worse-than-random 40% on the UCM dataset respectively (See Fig. 3).

Much better results were obtained when we compared our optimized pipeline to pipelines that use only a subset of our features: namely, the ones that use only features derived from sequence statistics and exclude the ones derived from learning PFSAs (recall that PFSAs are special HMMs we learn using our novel algorithms) from the disease categories as described in Methods in the main text, or using only the PFSA-based SLD features, or using simply the density of diagnostic codes (See Fig. 5, panel D). In all these cases we analyzed, our pipeline has a clearly demonstrable advantage (See Fig. 5, panel D) that is stable across databases, under reductions in sample sizes, and in balanced resampling experiments (See Fig. 5, panel C).

While it is difficult to explain the exact source of a modeling framework's performance, and even more difficult to explain non-performance, we can point to the following advantages that our approach has over existing techniques:

1. **Purely Classification Algorithms With No Pre-processing Do not Do well.** Pure classifiers such as random forests, gradient boosters, etc. are not time series modeling frameworks, and might not capture stochastic temporal patterns well. While features are not certainly assumed to be independent in these algorithms, it is problematic to learn patterns that do not appear at fixed time points in the diagnostic history.
2. **Lower Sample Complexity Compared to Deep Learning Frameworks.** Compared to LSTMs and RNNs, we are able to capture stochastic behavior with more compact models, which results in better sample complexity. In other words, if we have less data, our models do better, because we estimate fewer parameters.
3. **Designed Bottom-up for Learning Stochastic Processes.** It is easily demonstrated that LSTMs and RNNs, while good models of complicated time series in many cases, do not work well for data that are generated by stochastic processes, *i.e.* are sample paths of a hidden process.
4. **We May Have Missed Some Clever Transformation.** It is possible that extensive tuning or feature selection with LSTMs, RNNs or CNNs or some combination thereof, can replicate our performance, or even do better. There will always be that possibility, notwithstanding how much effort we put in to evaluate competing techniques. The authors welcome *future work in this direction that surpasses our performance reported here; this is only going to help the patients which is what matters.*

4.1. Feature Subset Evaluations & Code Density As A Feature

With regards to Fig. 5, panel D, we note that the PFSA based features by themselves are comparable to those engineered manually from sequence statistics (the latter include features such as the proportion of codes in a patient's history corresponding to specific disease categories, mean and variance of adjacent empty weeks *etc.*, see main text Table 3 in the main text for details), but the combined runs produce significantly superior results. Also, it is interesting to note that simply using the density of diagnostic codes in a child's history is quite predictive of future ASD diagnosis, with the AUC from using just the density of codes as a feature rising to over 75% in the Truven database at 150 weeks. However, it does not have stable predictive performance across databases, and is also the least performing predictor. We did not include code density in our combined feature set, since it has no effect once the rest of the features are combined.

5. Threshold Selection on ROC Curve

Once the ROC curve has been computed, we must choose a decision threshold to trade-off true positive rate and false positive rate. 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:

$$\text{accuracy} = \frac{t_p + t_n}{t_p + f_p + f_n + t_n} \quad (1)$$

$$\text{F1} = \frac{2t_p}{2t_p + f_p + f_n} \quad (2)$$

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 probably the one corresponding to 95% specificity, which is highlighted in Fig. 2C in the main text.

6. Note on Receiver 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 sake 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 1. *True positive rate, true negative rate, false positive rate, positive predictive value (PPV), and prevalence (ρ) are defined as:*

$$\text{TPR} = \frac{t_p}{P} = \frac{t_p}{t_p + f_n} \quad (3)$$

$$\text{TNR} = \frac{t_n}{N} = \frac{t_n}{t_n + f_p} \quad (4)$$

$$\text{FPR} = 1 - \text{TNR} \quad (5)$$

$$\text{PPV} = \frac{t_p}{t_p + f_p} \quad (6)$$

$$\rho = \frac{P}{N + P} \quad (7)$$

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

Note that TPR is also referred to as **recall** or **sensitivity**, and PPV is also referred to as **precision**. True negative rate is also known as **specificity**.

A **precision-recall curve**, or a PPV-sensitivity curve is a plot between PPV and TPR.

Denoting sensitivity by s , and specificity by c , it follows that:

$$\text{PPV} = \frac{t_p/P}{t_p/P + (f_p/N)(N/P)} = \frac{\text{TPR}}{\text{TPR} + ((N - t_n)/N)(N/P)} \quad (8)$$

$$\Rightarrow \text{PPV} = \frac{s}{s + (1 - c)(\frac{1}{\rho} - 1)} \quad (9)$$

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, if specificity and sensitivity are held constant. Also, if prevalence is limited to 2%, and specificity is held at 95%, then the maximum PPV is limited to:

$$\text{PPV} = s/(s + 2.45) \leq 1/3.45 \sim 29\% \quad (10)$$

This shows that for ASD screening, we can hope for a maximum PPV of ~29% at 95% specificity, if the prevalence is stable at around 2%.

Compare this with the PPV of 15.8% (M) and 18.8% (F) that we achieve at 51.8% sensitivity, where the specificity is held at 95% in Fig. 2C in the main text. Note here that M-Chat/F with follow-up has a PPV of 14.6% as reported by the recent CHOP study¹.

7. 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:

$$\text{TPR} = \frac{t_p}{P} = \frac{t'_p}{P'} \quad (11)$$

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 $\text{TPR} = 0$, 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.

We do have significant class imbalance in our datasets. This arises naturally from the low prevalence rate of ASD (small in the sense of comparison of sizes of the control and the positive cohorts). Thus, we validated if the performance of our predictive pipeline remains unchanged by replacing the full control cohort with a random sample of size equal to that of the positive cohort. The results, shown in Fig. 5C, illustrate that class imbalance has no appreciable effect on our pipeline, as far as the AUC metric is considered.

The precision-recall curves do get affected by class imbalance, or the prevalence, as shown by Eq (9). 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.

We have used this to our advantage. Note that since the PPV is affected by prevalence, a stratification of the total population with different prevalence in each sub-population suggests the possibility of a conditional choice of the operating point, thus boosting the overall PPV. We describe this approach in the sequel, in Section 9.1. First, we establish that our pipeline does not suffer from some important pitfalls arising in the workflows associated with ASD diagnosis, and how the diagnostic codes in Electronic Health Records (EHR) are generated.

8. Note on ASD Clinical Diagnosis & Uncertainty of EHR Record

With no precise laboratory test for ASD, most families experience the following sequence of events^{6–8}: 1) routine screening at 18 and 24 months of age identifies high risk, and is followed by 2) a diagnostic evaluation. The American Academy of Pediatrics (AAP) recommends screening all children for symptoms of ASD at 18 and 24 months of age in their primary care visits^{9,10}. However, results of a screening test are not diagnostic (*and hence do not produce an EHR diagnostic code*); they help the primary care provider identify children who are at risk for a diagnosis of ASD and require additional evaluation. The M-CHAT/F is the most studied and widely used tool for screening toddlers for ASD^{8,11}.

Unfortunately, children with milder symptoms are harder to screen for. The AAP warns that children with milder symptoms and/or average or above-average intelligence may not be identified with symptoms until school age, when differences in social language or personal rigidities affect function⁸.

8.1. Diagnostic Evaluations

Once a child is determined to be at risk for a diagnosis of ASD, either by screening or surveillance, a timely referral is needed for clinical diagnostic evaluation⁷, which will, on positive identification, assign a clinical diagnosis, and produce an EHR record.

The history of symptoms of ASD presentation in individual patients may be elucidated by questionnaires such as the Social Communication Questionnaire (SCQ), or Social Responsiveness Scale (SRS), or the Autism Diagnostic Interview-Revised (ADI-R)⁸. These questionnaires alone are insufficient for making a clinical diagnosis, but provide a structured approach to elicit symptoms. Validated observation tools used to provide structured data to

confirm a clinical diagnosis include the Autism Diagnostic Observation Schedule, Second Edition (ADOS-2)¹² and the Childhood Autism Rating Scale, Second Edition (CARS-2)¹³. Current guidance from the American Academy of Pediatrics⁸ notes that no single observation tool is universally appropriate, and that such tools are meant to support the application of the diagnostic criteria informed by history and other data.

At present, the Autism Diagnostic Interview-Revised (ADI-R) and Autism Diagnostic Observation Schedule (ADOS) are considered the “gold standard” tools to enable the diagnosis of ASD¹⁴. The true “gold standard” classification and diagnosis of autism is historically taken to be a multi-disciplinary team (MDT) clinical assessment, including use of the ADOS and ADI-R, as well as other assessments with consensus clinical judgment¹⁴. The MDT clinical diagnosis correct classification rate for ASD is approximately 80.8%. Thus, any individual tool that correctly classifies ASD at a rate of 80 % or over could be considered to be just as accurate as the “gold standard”¹⁴. With ADOS-2 and associated tools verifiably reaching this classification rate, the current APA guidance suggests that individual general pediatricians might hand out initial diagnoses if they are familiar with the relevant DSM diagnostic criteria. This simultaneously raises the prevalence, and the possibility that some diagnostic codes pertaining to ASD in medical history databases could be arising from less restrictive workflows, and thus might carry more uncertainty.

In our study, we checked if restricting the treatment cohort to children with at least two distinct ASD diagnostic codes in their medical histories instead of one (which significantly reduces the possibility of erroneous coding) changes the performance of the algorithm. The results shown in Fig. 5B illustrate that we have very little change in out-of-sample predictive performance, thus alleviating this concern.

8.2. Change In Diagnostic Criteria for ASD, Inclusion of PDD, Asperger, and Disambiguation From Unrelated Psychiatric Phenotypes

The DSM-5 established a single category of ASD to replace the subtypes of autistic disorder, Asperger syndrome, and pervasive developmental disorder not otherwise specified in the Diagnostic and Statistical Manual of Mental Disorders, Fourth Edition, Text Revision (DSM-IV-TR)⁸. This justifies our use of diagnostic codes from ICD9 299.X as specification of an ASD diagnosis, and use of GEMS mapping to 299.X for ICD10 codes when we encounter them. Future renditions of our pipeline will use purely ICD-10 specification, which does not change the algorithm, but merely how we input data into it.

It is interesting to note that we would be actually unable to discriminate between those phenotypes effectively for high predictability even if we wanted: in our initial efforts, we found it is very difficult to design a high performing pipeline that recognizes these sub-types separately.

The question then arises as to how well we can discriminate between ASD and other unrelated psychiatric phenotypes. Does our pipeline pick up on any psychiatric conditions, or is it specific to ASD? We directly evaluated this, by restricting the test control cohort to patients with at least one psychiatric code other than ASD. We get very high discrimination reaching AUCs over 90% at 100-125 weeks of age, which establishes that our pipeline is indeed largely specific to ASD.

8.3. Performance Comparison With M-CHAT/F

The M-CHAT/F is the most studied and widely used tool for screening toddlers for ASD^{8,11}.

Guthrie *et al.*¹ from the Children’s Hospital of Philadelphia (CHOP) demonstrate that when applied as a universal screening tool, M-CHAT/F has a sensitivity of 38.8%, specificity of 94.9% and PPV of 14.6%. This work is the only large-scale study of M-CHAT/F (n=20,375) we are aware of with sufficient follow-up after the age of four years to provide a reasonable degree of confidence in the sensitivity of M-CHAT/F.

Comparing the performance metrics achieved at different age groups across data sets and sexes for our pipeline (See main text Table 4 in the main text), we conclude that our approach produces a strictly superior PPV (exceeding M-CHAT/F PPV by at 14% (14.1-33.6%) when sensitivity and specificity are held at comparable values around the age of 26 months (\approx 112 weeks). We cannot compare at other operating points due to a lack of M-CHAT/F performance characterization anywhere else.

Apart from standalone performance, our proposed approach has several key advantages: it is clearly immune to parental educational level, and language barriers. Since access to insurance and medical records do get impacted by socio-economic variables, there is the possibility of some indirect impact from the demographic

makeup of the training datasets. But overall, diagnostic histories are free from biases that have historically plagued questionnaire-based screens⁸. Additionally, while M-CHAT/F is relatively easy and quick to administer, the issue of time and resource commitment cannot be ignored⁸. These factors conspire to produce reduced coverage, which in turn casts doubt upon the necessity of universal screening programs despite clear guidance on the contrary from the AAP¹.

Additionally, being functionally independent of the M-CHAT/F, we can take advantage of any population stratification induced by the M-CHAT/F results to significantly boost combined screening performance.

9. Improving Wait-times For Diagnostic Evaluations by Reducing False Positives in Routine Screening

While children with ASD can be diagnosed as toddlers^{15,16} (developmental concerns may show up before the first birthday^{17,18}), the mean age of diagnosis is over 4 years¹⁹. Since a clinical diagnosis of ASD requires the multi-step process described in the previous section, this delay mainly arises from extended wait-times and queues, which ultimately delays entry into early intervention (EI) programs. While time-consuming evaluations²⁰, cost of care²¹, lack of providers²², lack of comfort in diagnosing by primary care providers²², and other challenges, are all responsible to varying degrees that culminate in these delays⁶, one rather obvious source is the limited PPV of screening tests that are available today. With the PPV of M-CHAT/F being around 14.6%, over 85 out of 100 people flagged for diagnostic evaluation are false positives, leading to wait times that currently range from 3 months to 1 year. To make matters worse, access to care and resources are sparse except near urban centers. For example, only 7% of developmental pediatricians practice in rural areas, and some states do not even have a developmental pediatrician^{6,23}.

A key contribution of this work is to be able to significantly reduce the number of false positives without sacrificing specificity, and thus significantly improving wait-times and patient outcomes.

9.1. 4D Decision Optimization Using M-CHAT/F Population Stratification To Boost PPV

Assume that there are m sub-populations such that: the total number of positives and negatives, and the prevalences in each sub-population are given by P_i , N_i and ρ_i respectively, with $i \in \{1, \dots, m\}$. Let β_i be the relative size of the sub-populations. Thus, we have:

$$P = \sum_i P_i \quad (12)$$

$$N = \sum_i N_i \quad (13)$$

$$\beta_i = \frac{N_i + P_i}{N + P} \quad (14)$$

$$\rho_i = \frac{P_i}{N_i + P_i} = \frac{P_i}{\beta_i(N + P)} \quad (15)$$

Therefore, denoting the sensitivity and specificity of the sub-populations as s_i and c_i respectively, we have:

$$s = t_p/P = \frac{\sum_i t_p|_i}{P} = \frac{\sum_i (t_p|_i/P_i) \times (\beta_i \rho_i (P + N))}{P} = \sum_i s_i \beta_i \frac{\rho_i}{\rho} \quad (16)$$

Thus, we end up with:

$$s = \sum_{i=1}^m s_i \gamma_i \quad (17a)$$

$$c = \sum_{i=1}^m c_i \gamma'_i \quad (17b)$$

$$PPV = \frac{s}{s + (1 - c)(\frac{1}{\rho} - 1)} \quad (17c)$$

where we have denoted:

$$\gamma_i = \beta_i \frac{\rho_i}{\rho}, \text{ and } \gamma'_i = \beta_i \frac{1 - \rho_i}{1 - \rho} \quad (17d)$$

Now, using Table 2, we can compute the values for γ_i, γ'_i , as shown below.

Using the prevalence and stratification parameters calculated from the CHOP study (See main text Table 3¹), we can compute a conditional choice of sensitivity and specificity for our tool, in each sub-population to ultimately yield an overall performance significantly superior to M-CHAT/F. We carry out a four-dimensional search at the age the CHOP population stratification is reported (26 months or 112 weeks approximately) to identify the feasible region with $PPV > 14.6\%$, or sensitivity $> 38.8\%$ while keeping specificity $> 94.9\%$ where each of these dimensions represent the independent choice of sensitivity in the corresponding sub-population. For each set of 4 choices, the corresponding specificities are read-off from our computed ROC curve, and then the overall sensitivity, specificity and PPV are calculated using Eq. (17). The results are shown in Fig. 4, where we include the computations at 75 weeks, 125 weeks, and 150 weeks, with the same population stratification (although understandably the stratification will deviate from the values obtained at 26 months for those other ages).

An important assumption here is that the two tests are independent. Since M-CHAT/F is based on the detection of behavioral signals of developmental delay associated with autism via questionnaires completed by the primary care-givers, while our pipeline is based on physical comorbidities, independence is reasonable. Hence, we can simulate the application of the pipeline to each sub-population, and compute the overall performance quantities using a pre-computed ROC curve. Here we use the curve corresponding to the age in weeks, but average the male and female ROC curves, which are close as shown in Fig. 2 in the main text. The male-female averaging is necessary since the results from the CHOP study does not report sex stratified data.

We show the feasible region obtained by this computation in Fig. 4 of this document, and in main text Fig. 4 of the main text. Particularly, note that we get a PPV close to or higher than 30% at the high precision (HP) operating point, or a sensitivity above 55% for the high recall (HR) operating point, when we restrict specificities to above 95%.

It is important to note that Eq. (17) and hence the results are dependent on the population prevalence ρ . We report the dependence of the solution to the 4D optimization for population prevalence between 1.7% (CDC estimate⁸), and 2.23% (CHOP estimate¹). In particular, it is illuminating to compare these results directly with M-CHAT/F performance, as shown in Fig. 4, panels B and C in the main text. In panel C, we show that for any stable population prevalence between 1.7% and 2.24%, we can achieve nearly double the PPV without losing sensitivity, or increase the sensitivity by about 50% without sacrificing PPV, while holding not letting the specificity to drop below 94%.

10. Generating PFSA Models From Set of Input Streams with Variable Input Lengths

Our PFSA reconstruction algorithm² 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.

11. Probabilistic Finite State Automata Inference

11.1. Probabilistic Finite-State Automaton

Let Σ be a finite alphabet of symbols with size $|\Sigma|$. The set of sequences of length d over Σ is denoted by Σ^d . The set of finite but unbounded sequences over Σ is denoted by Σ^* , the Kleene star operation²⁴, i.e. $\Sigma^* = \bigcup_{d=0}^{\infty} \Sigma^d$. We use lower case Greek, for example σ or τ , for symbols in Σ , and lower case Latin, for example x or y , for sequences of symbols, i.e. $x = \sigma_1 \sigma_2 \dots \sigma_n$. We use $|x|$ to denote the length of x . The empty sequence is denoted by λ .

We denote the set of strictly infinite sequences over Σ by Σ^ω , and the set of strictly infinite sequences having x as prefix by $x\Sigma^\omega$. Let $S = \{x\Sigma^\omega : x \in \Sigma^*\} \cup \{\emptyset\}$, we can verify that S is a semiring²⁵ over Σ^ω . We use \mathcal{F} to denote the sigma algebra generated by S .

Definition 2 (Stochastic Process over Σ). A stochastic process over a finite alphabet Σ is a collection of Σ -valued random variables $\{X_t\}_{t \in \mathbb{N}}$ indexed by positive integers²⁶.

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

Definition 3 (Sequence-Induced Measure and Derivative). For a process \mathcal{P} , let $Pr_{\mathcal{P}}(x)$ or simply $Pr(x)$ denote the probability \mathcal{P} producing a sample path prefixed by x . The measure μ_x induced by a sequence $x \in \Sigma^*$ is the extension²⁵ to \mathcal{F} of the premeasure defined on the semiring \mathcal{S} given by

$$\forall x, y \in \Sigma^*, \mu_x(y\Sigma^\omega) \triangleq \frac{Pr(xy)}{Pr(x)}, \text{ if } Pr(x) > 0 \quad (18)$$

For any $d \in \mathbb{N}$, the d -th order derivative of a sequence x , written as ϕ_x^d , is defined to be the marginal distribution of μ_x on Σ^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 ϕ_x for short.

Definition 4 (Probabilistic Nerode Equivalence and Causal States²⁷). For any pair of sequences $x, y \in \Sigma^*$, 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**²⁸. 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 5 (Probabilistic Finite-State Automaton (PFSA)). A PFSA G is defined by a quadruple $(Q, \Sigma, \delta, \tilde{\pi})$, where Q is a finite set, Σ is a finite alphabet, $\delta : Q \times \Sigma \rightarrow \Sigma$ is called the transition map, and $\tilde{\pi} : Q \rightarrow \mathbf{P}_\Sigma$, where \mathbf{P}_Σ is the space of probability distributions over Σ , is called the transition probability. The entry of $\tilde{\pi}(q)$ indexed by σ is denoted by $\tilde{\pi}(q, \sigma)$.

Definition 6 (Transition and Observation Matrices). The transition matrix Π is the $|Q| \times |Q|$ matrix with the entry indexed by q, q' , written as $\pi_{q,q'}$, satisfying

$$\pi_{q,q'} \triangleq \sum_{\{\sigma \in \Sigma | \delta(q, \sigma) = q'\}} \tilde{\pi}(q, \sigma) \quad (19)$$

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

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

Definition 7 (Extension of δ and $\tilde{\pi}$ to Σ^*). For any $x = \sigma_1 \dots \sigma_k$, $\delta(q, x)$ is defined recursively by

$$\delta(q, x) \triangleq \delta(\delta(q, \sigma_1 \dots \sigma_{k-1}), \sigma_k) \quad (20)$$

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

$$\tilde{\pi}(q, x) \triangleq \prod_{i=1}^k \tilde{\pi}(\delta(q, \sigma_1 \dots \sigma_{i-1}), \sigma_i) \quad (21)$$

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

Definition 8 (Strongly Connected PFSA). We say a PFSA is strongly connected if the underlying directed graph is strongly connected²⁹. More precisely, a PFSA $G = (Q, \Sigma, \delta, \tilde{\pi})$ is strongly connected if for any pair of distinct states q and $q' \in Q$, there is an $x \in \Sigma^*$ 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 $\mathbf{v}^T \Pi = \mathbf{v}^T$. This is the **stationary distribution**^{30,31} of G and is denoted as φ_G , or φ if G is understood.

Definition 9 (Γ -Expression). We can encode the information contained in δ and $\tilde{\pi}$ by a set of $|Q| \times |Q|$ matrices $\Gamma = \{\Gamma_\sigma | \sigma \in \Sigma\}$, where

$$\Gamma_\sigma|_{q,q'} \triangleq \begin{cases} \tilde{\pi}(q, \sigma) & \text{if } \delta(q, \sigma) = q', \\ 0 & \text{if otherwise.} \end{cases} \quad (22)$$

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

Definition 10 (Sequence-Induced Distribution on States). For a PFSA $G = (Q, \Sigma, \delta, \tilde{\pi})$ and a distribution ρ_0 on Q , the **distribution on Q induced by a sequence x** is given by $\rho_{G, \rho_0}^T(x) = [\rho_0^T \Gamma_x]$ with $\rho_{G, \rho_0}(\lambda) = \rho_0$. The entry indexed by $q \in Q$ of the vector $\rho_{G, \rho_0}(x)$ is written as $\rho_{G, \rho_0}(x, q)$. When $\rho_0 = \rho_G$, the stationary distribution of G , we write $\rho_{G, \rho_0}(x)$ as $\rho_G(x)$, or simply as $\rho(x)$, if G is understood.

Definition 11 (Stochastic Process Generated by a PFSA). Let $G = (Q, \Sigma, \delta, \tilde{\pi})$ be a PFSA and let ρ_0 be a distribution on Q , the Σ -valued stochastic process $\{X_t\}_{t \in \Sigma}$ generated by G and ρ_0 satisfies that X_1 follows the distribution ρ_0 and X_{t+1} follows the distribution $\rho_{G, \rho_0}(X_1 \cdots X_t)$ for $t \in \mathbb{N}$.

For the rest of this paper, we will assume $\rho_0 = \rho_G$ if not specified otherwise. We can show that, when initialized with ρ_G , the process generated by a PFSA G is stationary and ergodic. We also note the, for the process generate by G , we have $\phi_x = \rho_G(x)^T \tilde{\Pi}$. Since $\rho_G(\lambda) = \rho_G$, the symbolic derivative of the empty sequence ϕ_λ is the stationary distribution on the symbols.

Definition 12 (Synchronizable PFSA and Synchronizing Sequence). A **synchronizing sequence** is a finite sequence that sends an arbitrary state of the PFSA to a fixed state³². To be more precise, let $G = (Q, \Sigma, \delta, \tilde{\pi})$ be a PFSA, we say a sequence $x \in \Sigma^*$ 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 13 (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 14. Consider a PFSA G over state set Q . For a give $\varepsilon > 0$, we say a sequence x is a ε -synchronizing sequence to a state $q \in Q$ if

$$\|\rho_G(x) - e_q\|_\infty \leq \varepsilon. \quad (23)$$

While there exists PFSA that is not synchronizable, we can show that an irreducible PFSA always has an ε -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 ε -synchronizing. These two are the underpinning properties for the inference algorithm of PFSA (See Alg. 3), because they imply that ϕ_x can be used to approximate $\tilde{\pi}(q)$ if x are properly prefixed and long enough.

Definition 15 (Joint ε -Synchronizing Sequence). Let G and H be two PFSA over state sets Q_G and Q_H , respectively. For a fixed ε , a sequence x is said to be **jointly ε -synchronizing** to $(q, r) \in Q_G \times Q_H$ if x is ε -synchronizing to q and to r simultaneously. We define

$$\Sigma_{\varepsilon, (q, r)}^d \triangleq \{x \in \Sigma^d : x \text{ jointly } \varepsilon\text{-synchronizing to } (q, r)\} \quad (24)$$

Definition 16 (Joint Pair of States). Let G and H be two PFSA over state sets Q_G and Q_H , respectively. Define

$$p_G(q, r) \triangleq \lim_{d \rightarrow \infty} p_G\left(\Sigma_{\varepsilon, (q, r)}^d\right) \quad (25)$$

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

$$Q_c \triangleq \{(q, r) \in Q_G \times Q_H : (q, r) \text{ is a } G\text{-joint pair}\} \quad (26)$$

The inference algorithm for PFSA is called **GenESS** for Generator Extraction Using Self-similar Semantics. With an input sequence x and a hyperparameter ε , **GenESS** 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. 3 for detail.

12. Sequence Likelihood Defect

Definition 17 (Entropy Rate and KL Divergence). By entropy rate of a PFSA, we mean the entropy rate of the stochastic process generated by the PFSA³³. Similarly, by KL divergence of two PFSA, we mean the KL divergence between the two processes generated by them³⁴. More precisely, we have

$$\mathcal{H}(G) = - \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{x \in \Sigma^d} p(x) \log p(x) \quad (27)$$

Algorithm 3: GenESeSS

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

Result: State set Q , transition map δ , and transition probability $\tilde{\pi}$

```

/* Step One: Approximate ε-synchronizing sequence */
```

- 1 Let $L = \lceil \log_{|\Sigma|} 1/\varepsilon \rceil$;
- 2 Calculate the **derivative heap** $\mathcal{D}_\varepsilon^x$ equaling $\left\{ \hat{\phi}_y^x : y \text{ is a sub-sequence of } x \text{ with } |y| \leq L \right\}$;
- 3 Let \mathcal{C} be the convex hull of $\mathcal{D}_\varepsilon^x$;
- 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** $x_{q_0}^{\text{id}} = x_0$ and the probability vector $d_{q_0} = \hat{\phi}_{x_0}^x$;
- 7 Let \tilde{Q} be the set of states that are just added and initialize it to be Q ;
- 8 **while** $\tilde{Q} \neq \emptyset$ **do**
- 9 Let $Q_{\text{new}} = \emptyset$ be the set of new states;
- 10 **for** $(q, \sigma) \in \tilde{Q} \times \Sigma$ **do**
- 11 Let $x = x_q^{\text{id}}$ 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 $Q_{\text{new}} = Q_{\text{new}} \cup \{q_{\text{new}}\}$ and $Q = Q \cup \{q_{\text{new}}\}$;
- 16 Associate to q_{new} the sequence identifier $x_{q_{\text{new}}}^{\text{id}} = x\sigma$ and the probability vector $d_{q_{\text{new}}} = d$;
- 17 Let $\delta(q, \sigma) = q_{\text{new}}$;
- 18 Let $\tilde{Q} = Q_{\text{new}}$;
- 19 Take a strongly connected subgraph of the labeled directed graph defined by Q and δ , 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 $\tilde{\pi}(q) = \llbracket (N[q, \sigma])_{\sigma \in \Sigma} \rrbracket$;
- 26 **return** $Q, \delta, \tilde{\pi}$;

and the KL divergence

$$\mathcal{D}_{KL}(G \parallel H) = \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{x \in \Sigma^d} p_G(x) \log \frac{p_G(x)}{p_H(x)} \quad (28)$$

whenever the limits exist.

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

$$\mathcal{H}(G) = \sum_{q \in Q} p_G(q) \cdot h(\tilde{\pi}(q)) \quad (29)$$

where $h(v)$ is the based-2 entropy of the probability vector v .

Consider two PFSA $G = (Q_G, \Sigma, \delta_G, \tilde{\pi}_G)$ and $H = (Q_H, \Sigma, \delta_H, \tilde{\pi}_H)$ with μ_G being absolutely continuous with respect to μ_H . Let Q_c be the set of G -joint pairs of states, we have

$$\mathcal{D}_{KL}(G \parallel H) = \sum_{(q, r) \in Q_c} p_G(q, r) D_{KL}(\tilde{\pi}_G(q) \parallel \tilde{\pi}_H(r)) \quad (30)$$

Definition 18 (Log-likelihood). *Let $x \in \Sigma^d$, the log-likelihood³³ of a PFSA G generating x is given by*

$$L(x, G) = -\frac{1}{d} \log p_G(x) \quad (31)$$

The calculation of log-likelihood is detailed in Alg. 4.

Algorithm 4: Log-likelihood

Data: A PFSA $G = (\Sigma, Q, \delta, \tilde{\pi})$ and a sequence x over alphabet Σ
Result: Log-likelihood $L(x, G)$ of G generating x

- 1 Calculate the state transition matrix Π and observation $\tilde{\Pi}$;
- 2 Calculate the stationary distribution over states φ_G of G from Π ;
- 3 Calculate the stationary distribution of alphabet $\phi_\lambda^T = \varphi_G^T \tilde{\Pi}$;
- 4 Initialize p by φ_G and q by ϕ_λ ;
- 5 Let $L = 0$;
- 6 **for** i from 1 to $|x|$ **do**
- 7 Let σ be the i -th entry of x ;
- 8 Let $L = L - \log q|_\sigma$;
- 9 Let $p^T = [\![p^T \Gamma_\sigma]\!]$ where Γ_σ is defined in 9;
- 10 Let $q^T = p^T \tilde{\Pi}$;
- 11 **return** $L/|x|$;

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*

$$L(x, H) \rightarrow \mathcal{H}(G) + \mathcal{D}_{KL}(G \parallel H) \quad (32)$$

in probability as $d \rightarrow \infty$.

Proof. We first notice that

$$\sum_{x \in \Sigma^d} p_G(x) \log \frac{p_G(x)}{p_H(x)} = \sum_{x \in \Sigma^{d-1}} \sum_{\sigma \in \Sigma} p_G(x) \varphi_G(x) \tilde{\Pi}_G \Big|_\sigma \log \frac{p_G(x) \varphi_G(x) \tilde{\Pi}_G \Big|_\sigma}{p_H(x) \varphi_H(x) \tilde{\Pi}_H \Big|_\sigma} \quad (33)$$

$$= \sum_{x \in \Sigma^{d-1}} p_G(x) \log \frac{p_G(x)}{p_H(x)} + \underbrace{\sum_{x \in \Sigma^{d-1}} p_G(x) \sum_{\sigma \in \Sigma} \varphi_G(x) \tilde{\Pi}_G \Big|_\sigma \log \frac{\varphi_G(x) \tilde{\Pi}_G \Big|_\sigma}{\varphi_H(x) \tilde{\Pi}_H \Big|_\sigma}}_{D_d} \quad (34)$$

By induction, we have $\mathcal{D}_{KL}(G \parallel H) = \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=1}^d D_i$, and hence by Cesàro summation theorem³⁵, we have $\mathcal{D}_{KL}(G \parallel H) = \lim_{d \rightarrow \infty} D_d$. Let $x = \sigma_1 \sigma_2 \dots \sigma_n$ be a sequence generated by G . Let $x^{[i-1]}$ is the truncation of x at the $(i-1)$ -th symbols, we have

$$-\frac{1}{n} \sum_{i=1}^n \log \varphi_H(x^{[i-1]}) \tilde{\Pi}_H \Big|_{\sigma_i} = \underbrace{\frac{1}{n} \sum_{i=1}^n \log \frac{\varphi_G(x^{[i-1]}) \tilde{\Pi}_G \Big|_{\sigma_i}}{\varphi_H(x^{[i-1]}) \tilde{\Pi}_H \Big|_{\sigma_i}}}_{A_{x,n}} - \underbrace{\frac{1}{n} \sum_{i=1}^n \log \varphi_G(x^{[i-1]}) \tilde{\Pi}_G \Big|_{\sigma_i}}_{B_{x,n}} \quad (35)$$

Since the stochastic process G generates is ergodic, we have

$$\lim_{n \rightarrow \infty} A_{x,n} = \lim_{d \rightarrow \infty} D_d = \mathcal{D}_{KL}(G \parallel H) \quad (36)$$

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

- [1] Guthrie, W. et al. Accuracy of Autism Screening in a Large Pediatric Network. *Pediatrics* **144** (2019).
- [2] Chattopadhyay, I. & Lipson, H. Abductive learning of quantized stochastic processes with probabilistic finite automata. *Philos Trans A* **371**, 20110543 (2013).
- [3] Breiman, L. Random forests. *Mach. Learn.* **45**, 5–32 (2001).
- [4] Friedman, J. H. Stochastic gradient boosting. *Comput. Stat. Data Anal.* **38**, 367–378 (2002).
- [5] Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural Comput.* **9**, 1735–1780 (1997).
- [6] Gordon-Lipkin, E., Foster, J. & Peacock, G. Whittling down the wait time: exploring models to minimize the delay from initial concern to diagnosis and treatment of autism spectrum disorder. *Pediatric Clinics* **63**, 851–859 (2016).
- [7] Penner, M., Anagnostou, E. & Ungar, W. J. Practice patterns and determinants of wait time for autism spectrum disorder diagnosis in canada. *Molecular autism* **9**, 16 (2018).
- [8] Hyman, S. L., Levy, S. E., Myers, S. M. et al. Identification, evaluation, and management of children with autism spectrum disorder. *Pediatrics* **145** (2020).
- [9] Johnson, C. P., Myers, S. M. et al. Identification and evaluation of children with autism spectrum disorders. *Pediatrics* **120**, 1183–1215 (2007).

- [10] Zwaigenbaum, L. *et al.* Early intervention for children with autism spectrum disorder under 3 years of age: recommendations for practice and research. *Pediatrics* **136**, S60–S81 (2015).
- [11] Robins, D. L. *et al.* Validation of the modified checklist for autism in toddlers, revised with follow-up (m-chat-r/f). *Pediatrics* **133**, 37–45 (2014).
- [12] Esler, A. N. *et al.* The autism diagnostic observation schedule, toddler module: standardized severity scores. *Journal of Autism and Developmental Disorders* **45**, 2704–2720 (2015).
- [13] Chlebowski, C., Green, J. A., Barton, M. L. & Fein, D. Using the childhood autism rating scale to diagnose autism spectrum disorders. *Journal of autism and developmental disorders* **40**, 787–799 (2010).
- [14] Falkmer, T., Anderson, K., Falkmer, M. & Horlin, C. Diagnostic procedures in autism spectrum disorders: a systematic literature review. *European child & adolescent psychiatry* **22**, 329–340 (2013).
- [15] Lord, C. *et al.* Autism from 2 to 9 years of age. *Archives of general psychiatry* **63**, 694–701 (2006).
- [16] Kleinman, J. M. *et al.* Diagnostic stability in very young children with autism spectrum disorders. *Journal of autism and developmental disorders* **38**, 606–615 (2008).
- [17] Bolton, P. F., Golding, J., Emond, A. & Steer, C. D. Autism spectrum disorder and autistic traits in the avon longitudinal study of parents and children: precursors and early signs. *Journal of the American Academy of Child & Adolescent Psychiatry* **51**, 249–260 (2012).
- [18] Kozlowski, A. M., Matson, J. L., Horovitz, M., Worley, J. A. & Neal, D. Parents' first concerns of their child's development in toddlers with autism spectrum disorders. *Developmental neurorehabilitation* **14**, 72–78 (2011).
- [19] Baio, J. Prevalence of autism spectrum disorder among children aged 8 years-autism and developmental disabilities monitoring network, 11 sites, united states, 2010 (2014).
- [20] Kalb, L. G. *et al.* Determinants of appointment absenteeism at an outpatient pediatric autism clinic. *Journal of Developmental & Behavioral Pediatrics* **33**, 685–697 (2012).
- [21] Bisgaier, J., Levinson, D., Cutts, D. B. & Rhodes, K. V. Access to autism evaluation appointments with developmental-behavioral and neurodevelopmental subspecialists. *Archives of pediatrics & adolescent medicine* **165**, 673–674 (2011).
- [22] Fenikilé, T. S., Ellerbeck, K., Filippi, M. K. & Daley, C. M. Barriers to autism screening in family medicine practice: a qualitative study. *Primary health care research & development* **16**, 356–366 (2015).
- [23] Althouse, L. A. & Stockman, J. A. Pediatric workforce: A look at pediatric nephrology data from the american board of pediatrics. *The Journal of pediatrics* **148**, 575–576 (2006).
- [24] Hopcroft, J. E. *Introduction to automata theory, languages, and computation* (Pearson Education India, 2008).
- [25] Klenke, A. *Probability theory: a comprehensive course* (Springer Science & Business Media, 2013).
- [26] Doob, J. *Stochastic processes*. Wiley publications in statistics (Wiley, 1990). URL <https://books.google.com/books?id=7Bu8jgECAAJ>.
- [27] Chattopadhyay, I. & Ray, A. Structural transformations of probabilistic finite state machines. *International Journal of Control* **81**, 820–835 (2008).
- [28] Chattopadhyay, I. & Lipson, H. Data smashing: uncovering lurking order in data. *Journal of The Royal Society Interface* **11**, 20140826 (2014).
- [29] Bondy, J. & Murty, U. Graph theory (2008). *Grad. Texts in Math* (2008).
- [30] Vidyasagar, M. *Hidden markov processes: Theory and applications to biology*, vol. 44 (Princeton University Press, 2014).
- [31] Kai, L. C. *Markov Chains: With Stationary Transition Probabilities* (Springer-Verlag, 1967).
- [32] Trahtman, A. N. The road coloring and Černý conjecture. In *Proc. of Prague stringology conference*, vol. 1, 12 (Citeseer, 2008).
- [33] Cover, T. M. & Thomas, J. A. *Elements of information theory* (John Wiley & Sons, 2012).
- [34] 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).
- [35] Hardy, G. Divergent series, with a preface by J. E. Littlewood and a note by G. H. Hardy and J. E. Littlewood, reprint of the revised (1963) edition. *Éditions Jacques Gabay, Sceaux* (1992).