

# Supplementary Information:

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

Dmytro Onishchenko<sup>1</sup>, Yi Huang<sup>1</sup>, James van Horne<sup>1</sup>, Peter J. Smith<sup>4,7</sup>, Michael M. Msall<sup>5,6</sup> and Ishanu Chattopadhyay<sup>1,2,3★</sup>

<sup>1</sup>Department of Medicine, University of Chicago, Chicago, IL USA

<sup>2</sup>Committee on Genetics, Genomics & Systems Biology, University of Chicago, Chicago, IL USA

<sup>3</sup>Committee on Quantitative Methods in Social, Behavioral, and Health Sciences, University of Chicago, Chicago, IL USA

<sup>4</sup>Department of Pediatrics, Section of Developmental and Behavioral Pediatrics, University of Chicago, Chicago, IL USA

<sup>5</sup>Department of Pediatrics, Section Chief of Developmental and Behavioral Pediatrics, University of Chicago, Chicago, IL USA

<sup>6</sup>Joseph P. Kennedy Research Center on Intellectual and Neurodevelopmental Disabilities, University of Chicago, Chicago, IL USA

<sup>7</sup>Executive Committee Chair, American Academy of Pediatrics' Section on Developmental and Behavioral Pediatrics

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

## CONTENTS

<b>I</b>	<b>Detailed Mathematical Approach</b>	9
I-A	Time-series Modeling of Diagnostic History . . . . .	9
I-B	Step 1: Partitioning The Human Disease Spectrum . . . . .	9
I-C	Step 2: Model Inference & The Sequence Likelihood Defect . . . . .	10
I-D	Step 3: Risk Estimation Pipeline With Semi-supervised & Supervised Learning Modules .	12
<b>II</b>	<b>Comparison With State of the Art Off-the-shelf ML Algorithms</b>	12
<b>III</b>	<b>Pipeline Variations, Feature Subsets and Neural Network (NN) Post-processing</b>	12
III-A	Feature Subset Evaluations & Code Density As A Feature . . . . .	14
<b>IV</b>	<b>Threshold Selection on ROC Curve</b>	14
<b>V</b>	<b>Note on Receiver Operating Characteristics (ROC) and Precision-recall Curves</b>	15
<b>VI</b>	<b>Effect of Class Imbalance</b>	19
<b>VII</b>	<b>Note on ASD Clinical Diagnosis &amp; Uncertainty of EHR Record</b>	19
VII-A	Diagnostic Evaluations . . . . .	19
VII-B	Change In Diagnostic Criteria for ASD, Inclusion of PDD, Asperger, and Disambiguation From Unrelated Psychiatric Phenotypes . . . . .	20
VII-C	Performance Comparison With M-CHAT/F . . . . .	20

<b>VIII</b>	<b>Improving Wait-times For Diagnostic Evaluations by Reducing False Positives</b>	21
VIII-A	4D Decision Optimization Using M-CHAT/F Population Stratification To Boost PPV . . . . .	21
<b>IX</b>	<b>Generating PFSA Models From Set of Input Streams with Variable Input Lengths</b>	23
<b>X</b>	<b>Probabilistic Finite State Automata Inference</b>	23
X-A	Probabilistic Finite-State Automaton . . . . .	23
<b>XI</b>	<b>Sequence Likelihood Defect</b>	25
<b>XII</b>	<b>Pipeline Optimization</b>	27
XII-A	Input Data Format . . . . .	27
XII-B	Algorithms . . . . .	27
<b>XIII</b>	<b>Example Run with Released Application</b>	27
XIII-A	Prerequisites & Installation . . . . .	27
XIII-B	EHR data format . . . . .	28
XIII-C	Sample Python code risk estimation . . . . .	28
XIII-D	Sample Python script risk estimation . . . . .	28

#### LIST OF TABLES

I	Disease Categories With Detailed Set of ICD9 Codes Used . . . . .	3
II	Boosted Sensitivity, Specificity and PPV Achieved at <b>150 weeks</b> Conditioned on M-CHAT/F Scores . . . . .	22
III	Population Stratification Results on large M-CHAT/F Study(n=20,375) (14) . . . . .	22
IV	$\gamma, \gamma'$ Computed from Population Stratification Recorded In M-CHAT/F Study (14) ( $\rho = 0.0223$ ) . . . . .	22

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

Category	Description	Constituent ICD9 Codes
Hematologic	Diseases Of The Blood And Blood-Forming Organs	286.9 286.7 286.6 283.19 283.10 283.11 283.9 283 283.1 284.9 284.8 284.81 284.0 284.89 284.09 284 284.01 282.2 287.49 287.41 287.39 287.4 287.5 287.32 287.3 287.30 287.31 286.3 286.2 286.1 286.0 286.4 282.1 282.6 282.5 282.41 282.42 282.68 282.69 282.62 282.63 282.60 282.61 282.64 282 282.8 287.33 281.2 281.3 280.0 282.9 285.8 285.9 280.9 284.2 285.1 285.2 285.3 280.1 285.22 285.21 282.3 276.5 285 283.0 285.29 280.8 282.7 282.40 282.49 284.1 284.19 284.12 284.11 281.8 281.9 281.4 281.0 281.1 286.5 287 287.8 287.9 287.2 287.0 287.1 285 289.52 289.50 289.51 289.59 289.4 289.5 289.81 289.83 289.82 289.89 289 289.7 289.8 289.9
Psychiatric	Mental Disorders (Except ASD)	290 through 319 (except 299.x)
Metabolic	Metabolic Disorders (Distinct from respiratory, digestive and immunological conditions)	273.4 270 270.2 270.3 712.11 712.10 712.13 712.12 712.15 712.14 712.17 712.16 712.19 712.18 712.31 712.30 712.37 712.36 712.35 712.34 712.38 712.33 712.32 712.28 712.29 712.24 712.25 712.26 712.27 712.20 712.21 712.22 712.23 712.39 712.1 712.3 712.2 277.6 275.1 277.85 277.87 270.7 270.6 276.6 276.7 276.4 276.2 276.3 276.0 275.41 276.1 276.8 276.9 276.69 275.5 275.42 271.1 330.2 272.7 271 274.89 712.85 274.81 274.82 712.99 712.98 274.01 274.00 274.03 274.02 712.91 712.90 712.93 712.92 712.95 712.94 712.97 712.96 712.88 712.89 274.10 274.11 712.82 712.83 712.80 712.81 712.86 712.87 712.84 274.19 712.9 712.8 274.0 274.1 274.8 274.9 271.2 275.01 270.5 270.4 278.8 272.3 275.03 275.09 271.3 272.6 272.5 278.1 271.8 277.5 263.0 263.2 262 260 261 263 263.1 269.8 269.9 263.8 263.9 269 277.7 272.2 272.3 272.6 272.7 272.8 277.8 275.49 275.2 277.88 275.4 269.3 275.9 275.8 277.9 277.89 251.2 251.1 251.0 278.01 278.00 278.03 270.8 270.9 278.0 278.02 277.86 270.1 275.3 277.1 277.81 277.82 272 272.1 277.2 272.4 272.9 273.9 273.8 268.1 265.2 268.0 268.2 268 265.0 265.1 266.1 266.0 266.2 266.9 264.3 264.2 264.1 264.0 264.7 264.6 264.5 264.4 264.9 264.8 268.9 267 266 265 264 269.2 269.0 269.1 278.2 278.3 278.4
Cardiovascular	Diseases Of Arteries, Arterioles, And Capillaries	442.89 441.6 442.84 442.82 442.83 442.8 441.03 441.02 441.01 441.00 441 414.19 414.12 442 414.10 414.11 447.70 447.71 447.72 447.73 414.1 442.81 441.9 442.1 442.0 442.3 442.2 441.2 441.3 441.0 441.1 442.9 441.7 441.4 441.5 437.3 447.7 443.29 443.23 443.22 443.21 443.24 443.2 444.9 444.8 444.81 444.2 444.1 444.0 444.0 444.89 444 444.22 444.21 445.81 440.31 440.30 440.32 414.01 414.00 414.03 414.02 414.05 414.04 414.07 414.06 445.89 411.81 445.02 445.01 440.24 440.22 440.23 440.20 440.21 440 445 440.29 441.0 441.2 440.4 440.3 440.2 440.1 440.0 440.9 440.8 445.8 445.0 414.3 414.4 426.54 426.53 426.52 426.51 426.50 426.13 426.12 426.11 426.10 426.89 426.9 426.8 426.81 426.3 426.2 426.1 426.0 426.7 426.6 426.5 426.4 427.61 427.60 427.5 427.89 427.69 427.32 427.41 427.9 427.81 427.8 427.4 427.4 427.3 427.0 427.1 425.8 425.9 425.4 425.4 425.7 425.0 425.1 425.2 425.3 425.6 438.51 438.52 438.53 438.50 290.4 431.0 438.42 438.41 438.40 432.9 290.43 290.42 290.41 290.40 432.0 432.1 433.00 433.01 434.9 346.61 346.60 346.63 346.62 433.80 433.81 433.11 433.10 438.32 438.30 438.31 434.0 433.91 433.90 430.0 434.10 434.11 434.1 438.21 438.20 438.22 433.20 433.21 438.6 438.7 438.4 438.5 438.2 438.3 438.0 438.1 438.8 438.9 438.436 434.90 434 435 432 433 430 434.91 434.01 434.00 438.10 438.11 438.12 438.13 438.14 438.19 433.31 433.30 438.85 438.84 438.83 438.82 438.81 438.89 433.9 433.8 433.1 433.0 433.3 433.2 437.453.5 435.2 435.1 435.0 435.3 453.5 453.4 453.3 453.2 453.1 453.0 453.9 453.8 453.52 453.51 453.50 453.79 453.71 453.73 453.72 453.75 453.74 453.77 453.76 453.84 453.89 453.40 453.41 453.42 453.81 453.82 453.83 415.11 453.85 453.86 453.87 405.1 405.1 404.9 404.9 403.11 402.00 402.01 404.1 404.0 402.1 402.0 402.0 403.0 405.99 402.9 405.91 402.91 402.90 405.11 401.0 401.1 404.00 404.01 404.02 404.03 405.19 401.9 405 404.0 403 402 401 405.9 403.01 403.00 402.11 402.10 403.10 404.13 404.12 404.11 404.10 405.01 403.9 405.09 403.90 437.2 403.1 403.91 404.93 404.92 404.91 404.90 448 458 458.0 458.2 458.1 458.9 458.29 458.21 426.82 429.71 410.01 410.00 410.02 410.41 410.40 410.42 410.22 410.21 410.20 429.7 410.70 410.71 410.72 429.79 410.92 410.90 410.91 410.30 410.31 410.32 410.12 410.10 410.11 410.52 410.50 410.51 410.4 410.5 410.6 410.7 410.0 410.1 410.2 410.3 410.8 410.9 411.0 410.62 410.61 410.60 410 412 410.81 410.80 410.82 424.1 424.0 424.2 424.29.89 429.429.1 429.5 429.6 429.8 429.9 459.7 276.5 429.2 429.3 428.9 428.4 428.8 428.1 428.0 428.20 428.20 428.20 428.22 428.23 428.00 459.89 448.1 454 455 455.9 455.8 454.8 454.9 455.1 455.0 455.3 455.2 455.5 455.4 455.7 455.6 454.2 447 454.0 454.1 757.32 447.8 447.9 448.9 447.4 447.5 447.0 447.2 447.3 414 413.1 413.0 413.9 411.89 411.1 411 414.9 413 414.8 411.8 443.89 443 459.2 443.8 443.9 443.81 459.81 443.0 416.2 415.19 415.1 415.13 415.12 416.
Nutrition	Nutrition, metabolism, and development	7830,78321,7833,78340,78342,7837,7839

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



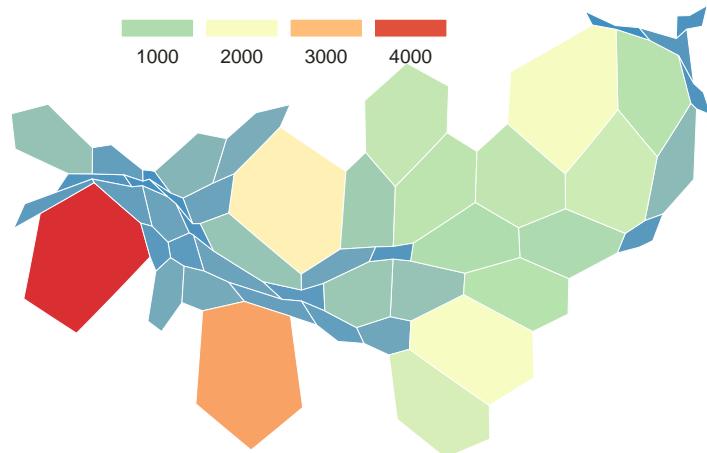




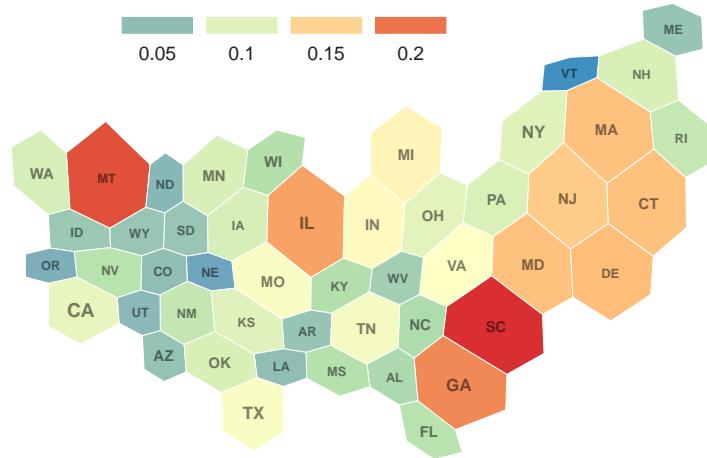




**A. Autism Insurance Claims 2003-2013**  
 (source: Truven Marketscan)



**B. Autism Prevalence in US (Population Normalized)**



SI-Fig. 1. **ASD Occurrence Patterns** Panel A illustrates the spatial distribution of ASD insurance claims, and panel B 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 (88).

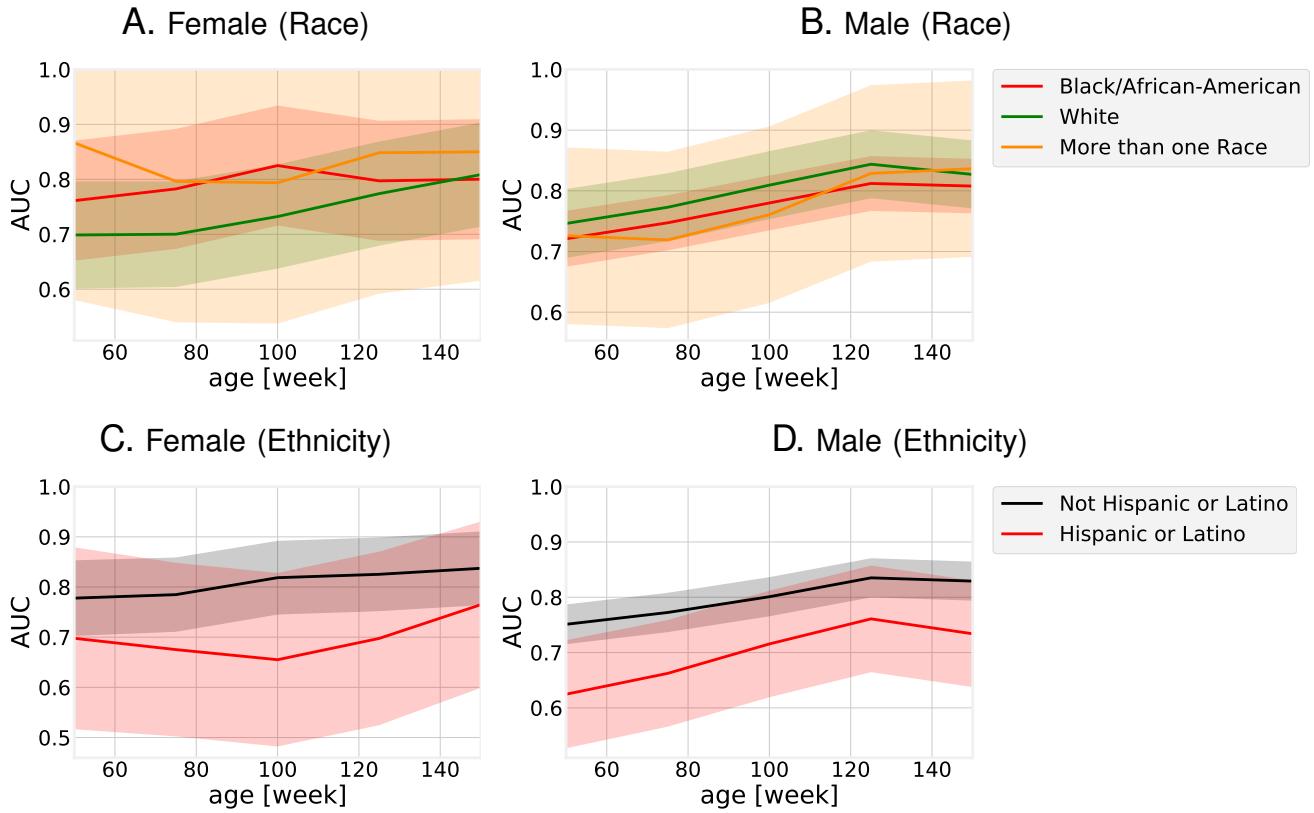
## I. DETAILED MATHEMATICAL APPROACH

### A. Time-series Modeling of Diagnostic History

Individual diagnostic histories can have long-term memory (95), 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) (31). 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 X). 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.

### B. Step 1: Partitioning The Human Disease Spectrum

We begin by partitioning the human disease spectrum into 17 non-overlapping categories. Each category is defined by a set of diagnostic codes from the International Classification of Diseases, Ninth Revision (ICD9)



SI-Fig. 2. **Effect of Race and Ethnicity on Predictive Performance with 95% Confidence Bounds from the UCM dataset.** Panels A and B show the variation of AUC achieved in out-of-sample data in three race-based population groups (White, African-American and multi-racial). We find no significant differences. Panels C and D show the performance in Hispanic vs non-Hispanic sub-populations. We find that children with Hispanic background have a lower AUC, but the differences are not significant. Other races/ethnicities were not considered due to lack of sufficient data.

(See Table SI-I 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) (71 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; we want our algorithm to seek out the important patterns without any manual curation of the input data. The limitation of the set of phenotypes to 9835 unique codes arises from excluding patients from the database who have very few and rare codes that will skew the statistical estimates. As shown in Table 1a in the main text, we exclude a very small number of patients, and who have very short diagnostic histories with a very small number of codes.

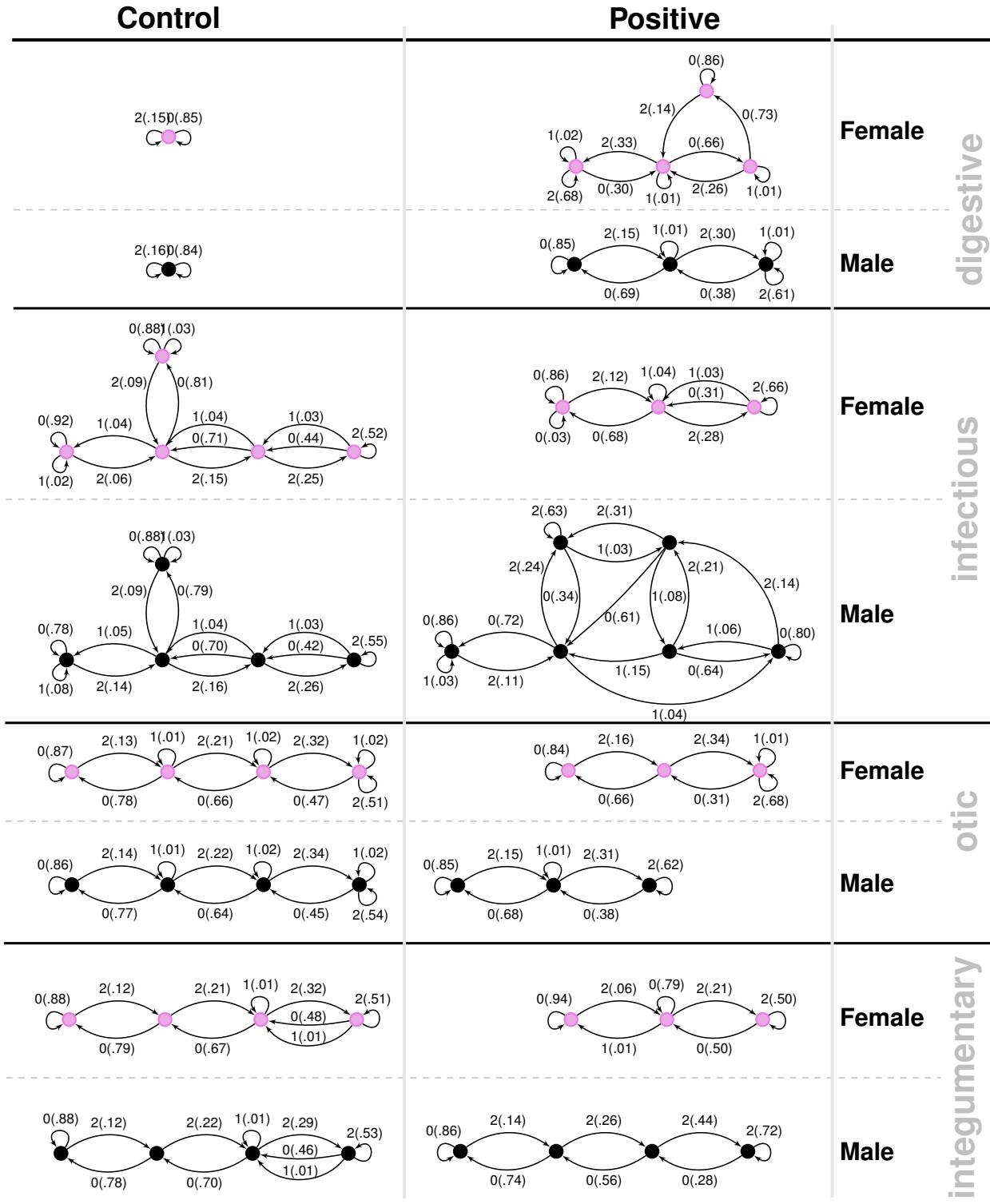
For each patient, the past medical history is a sequence  $(t_1, x_1), \dots, (t_m, x_m)$ , where  $t_i$  are timestamps and  $x_i$  are ICD9 codes diagnosed at time  $t_i$ . We map individual patient history to a three-alphabet categorical time series  $z^k$  corresponding to the disease category  $k$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

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

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

### C. Step 2: Model Inference & The Sequence Likelihood Defect

The mapped series, stratified by sex, disease-category, and ASD diagnosis-status are considered to be independent sample paths, and we want to explicitly model these systems as specialized HMMs (PFSAs). We



SI-Fig. 3. 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.

model the positive and the control cohorts for each sex, and in each disease category separately, ending up with a total of 68 HMMs at the population level (17 categories, 2 sexes, 2 cohort-types: positive and control, SI-Fig. 3 in the supplementary text provides some examples). Each of these inferred models is a PFSA; a directed graph with probability-weighted edges, and acts as an optimal generator of the stochastic process driving the sequential appearance of the three letters (as defined by Eq. (1)) corresponding to each sex, disease category, and cohort-type (See Section X in the Supplementary text for background on PFSA 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 (34, 35) between probability distributions to a divergence  $\mathcal{D}_{\text{KL}}(G||H)$  between ergodic stationary categorical stochastic processes (36)  $G, H$  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 X, Theorem 6 and 7):

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

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

where  $\mathcal{H}(\cdot)$  is the entropy rate of a process (34). Importantly, Eq. (4) shows that the computed likelihood has an additional non-negative contribution from the divergence term when we choose the incorrect generative process. Thus, if a patient is eventually going to be diagnosed with 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,  $\Delta^j$ ) 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. (5). The higher this likelihood defect, the higher the similarity of diagnosis history to that of children with autism.

#### D. 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 (84). The complete list of 165 features used is provided in Tab. Ib in the main text.

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. 1B in the main text shows the top 15 features ranked in order of their relative importance (relative loss in performance when dropped out of the analysis).

## II. COMPARISON WITH STATE OF THE ART OFF-THE-SHELF ML ALGORITHMS

### SeaGreen4

Off the shelf algorithms with little or no pre-processing, *i.e.*, using the diagnostic codes themselves are time-stamped categorical features failed to produce clinically relevant performance (See SI-Fig. 5). Classifiers such as random forests (75), and gradient boosters (84) 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 LSTM (86) 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.

## III. PIPELINE VARIATIONS, FEATURE SUBSETS AND NEURAL NETWORK (NN) POST-PROCESSING

In addition to the naive baseline approaches, we also evaluated the performance achievable with LSTMs (denoted as LSTMb in SI-Fig. 5) that use identical pre-processing as our pipeline, *i.e.*, representation of diagnostic histories

Model: "sequential\_4"

Layer (type)	Output Shape	Param #
lstm_1 (LSTM)	(None, 150)	180600
dense_5 (Dense)	(None, 32)	4832
dense_6 (Dense)	(None, 1)	33
Total params:	185,465	
Trainable params:	185,465	
Non-trainable params:	0	

SI-Fig. 4. The simplest LSTM investigated as a baseline. More complex models have significantly larger number of trainable parameters (1 to 10 Million). In contrast our pipeline has 13,744 trainable parameters.

as trinary sequences in 18 categories for each patient, and achieved ~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. 5).

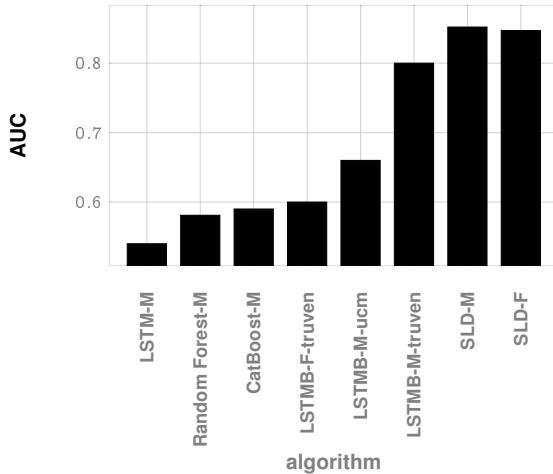
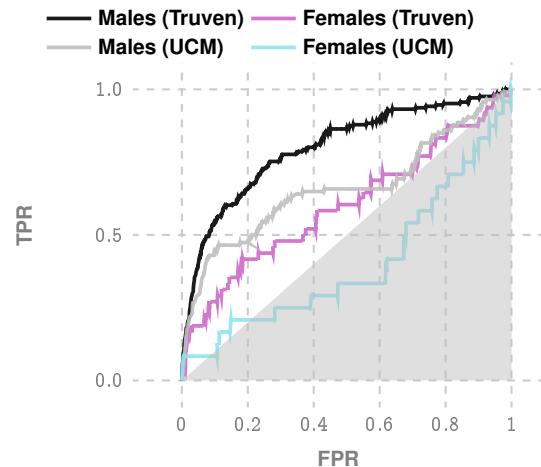
Much better results were obtained when we compared our optimized pipeline to pipelines that use only a subset of our features: namely, the ones that use only features derived from sequence statistics and exclude the ones derived from learning PFSA (recall that PFSA are special HMMs we learn using our 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. 6, panel D). In all these cases, our pipeline has a clearly demonstrable advantage (See Fig. 6, panel D) that is stable across databases, under reductions in sample sizes, and in balanced re-sampling experiments (See Fig. 6, panel C).

While it is difficult to explain the exact source of a modeling framework's 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.

Thus, there are two related issues: NN models generally have a relatively large number of nominal free parameters, requiring a substantial amount of data to train. In our case, we indeed have a large number of control examples, but a relatively smaller number of positive examples (children who get diagnosed with autism eventually). The amount of training data, even with our large patient database, is not enough to train anything but the simplest of the NN architectures, and even then the number of parameters for such NN models is substantially larger compared to our architecture. This is shown in our example (See SI-Fig. 4), where one of the LSTM models we investigate (the simplest one we tested) has 185,465 parameters, whereas our complete end to end pipeline has  $7025+6719=13,744$  trainable parameters (an order of magnitude less, where the first contribution is from inferred PFSA models, the second from the ensemble of decision trees in the final gradient boosting model). Additionally, the number of parameters for deep learning models rapidly increase to millions or tens of millions as the models get more complicated.

Also note that our framework is adaptive (See examples of PFSA models inferred in SI-Fig. 3), and the network structures of the inferred PFSA models is inferred from data, and not fixed a priori; it produces more states where it needs, thus adapting the model resolution to more complex contexts, and obtaining significantly more compact yet effective representations. We must however add the caveat that with sufficient experience it might be possible to devise a specific NN architecture with similar performance to our framework; however, in our

**A. Sample of Baseline Approaches with AUC > 0.5****B. ROC Curves for LSTMB (LSTM with pre-processing)**

SI-Fig. 5. 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.

approach such artful insight is unnecessary.

We also found that NN models generalize poorly across databases for the problem at hand; the performance is relatively better in the Truven dataset (See SI-Fig. 5), but worse in an independent database.

#### A. Feature Subset Evaluations & Code Density As A Feature

With regards to Fig. 6, 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 Ib 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.

#### IV. 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:

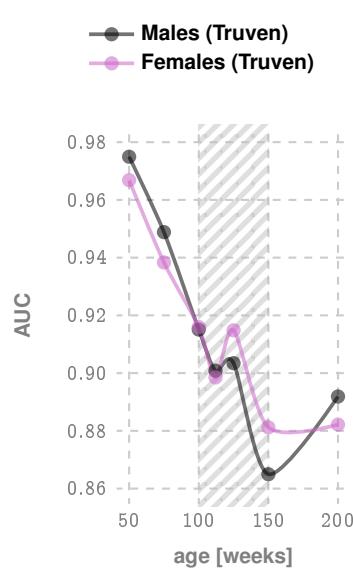
$$\text{accuracy} = \frac{t_p + t_n}{t_p + f_p + f_n + t_n} \quad (6)$$

$$F1 = \frac{2t_p}{2t_p + f_p + f_n} \quad (7)$$

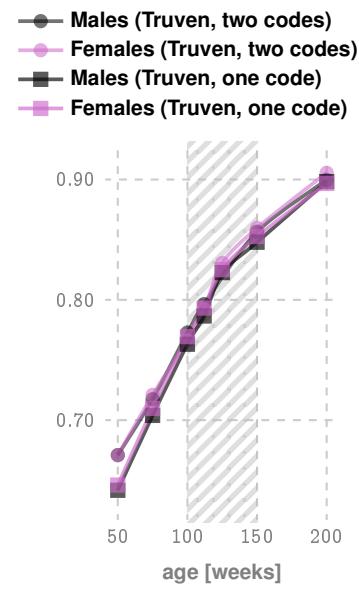
The F1-score is the same as accuracy where the number of true negatives is the same as the number of true positives, thus partially correcting for the class imbalance.

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. 1C in the main text.

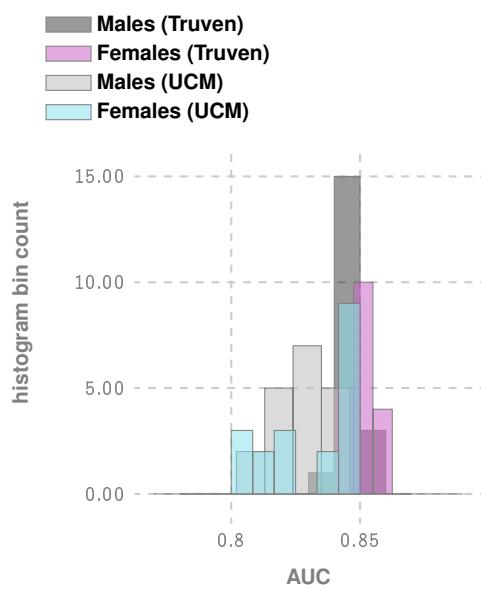
### A. Disambiguation of Autism Diagnosis from Other Psych. Phenotypes



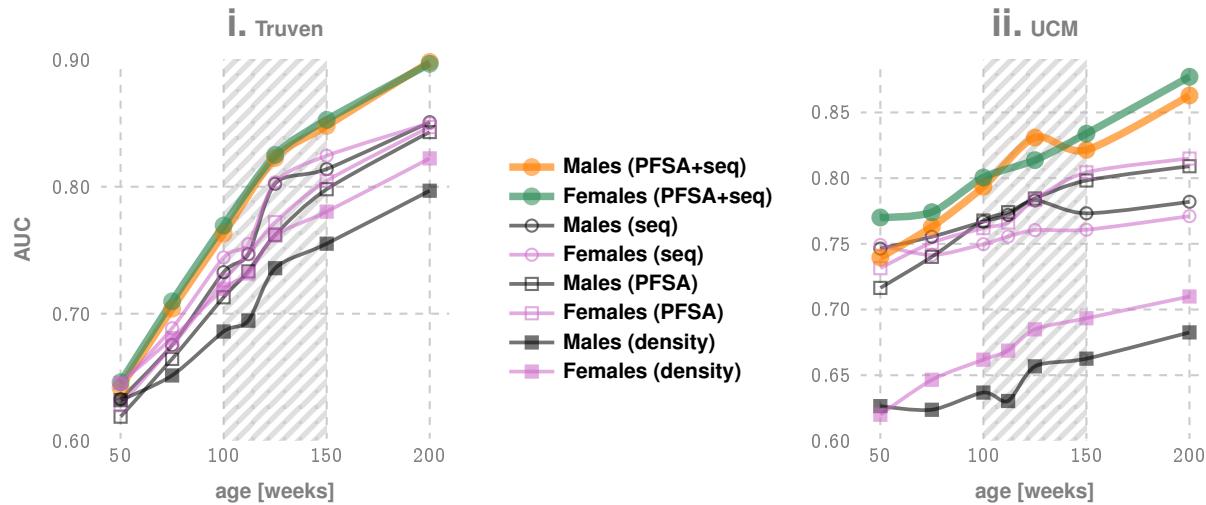
### B. Comparison of Performance with One vs Two ASD Diagnostic Codes



### C. AUC Distribution with Matched Control & Treatment Population Sizes



### 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. 6. **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.

### V. 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 same of completeness, we introduce the relevant definitions:

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

**Definition 1.** *True positive rate, true negative rate, false positive rate, positive predictive value (PPV), and*



SI-Fig. 7. **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).

**prevalence ( $\rho$ ) are defined as:**

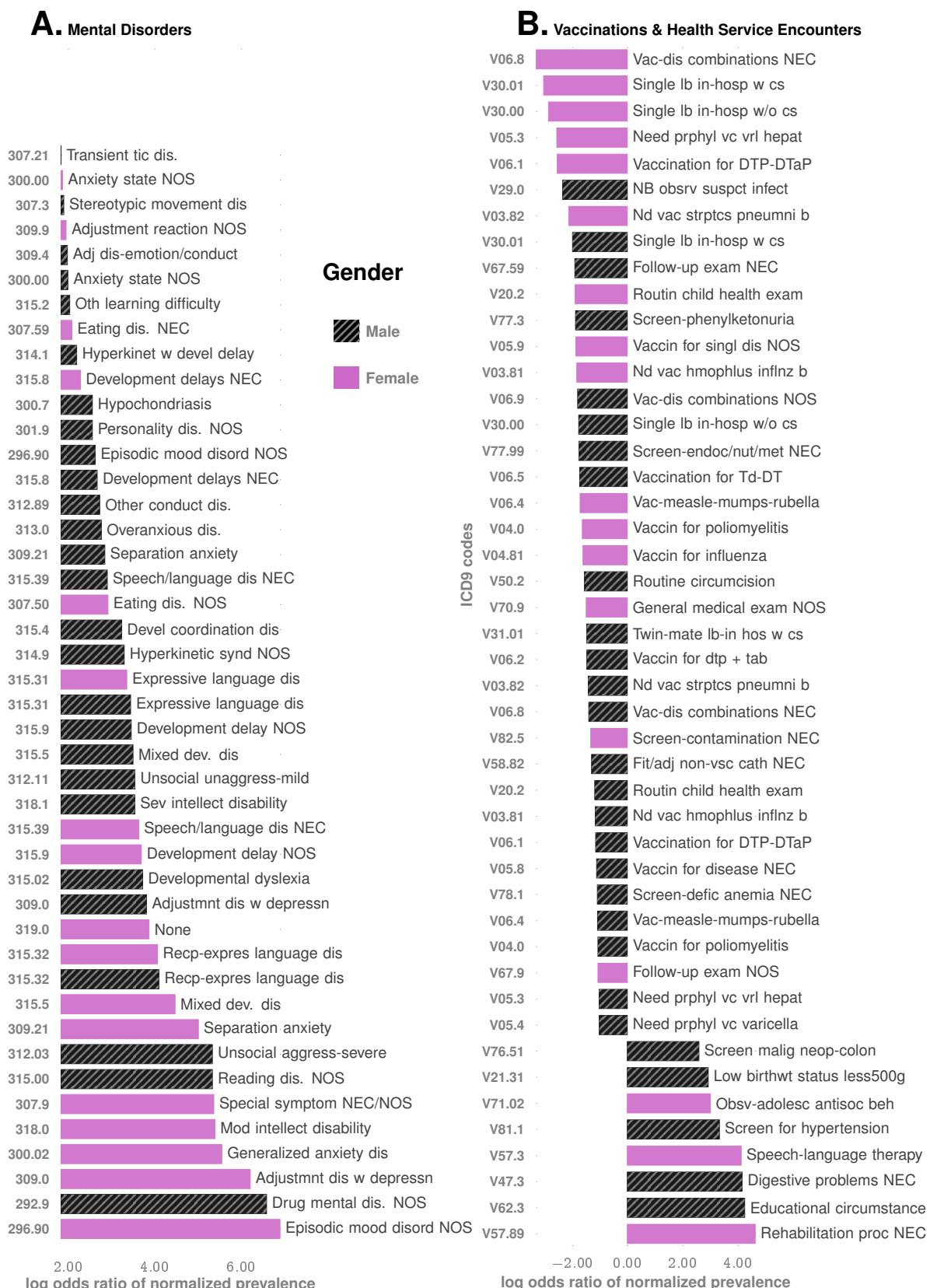
$$TPR = \frac{t_p}{P} = \frac{t_p}{t_p + f_n} \quad (8)$$

$$TNR = \frac{t_n}{N} = \frac{t_n}{t_n + f_p} \quad (9)$$

$$FPR = 1 - TNR \quad (10)$$

$$PPV = \frac{t_p}{t_p + f_p} \quad (11)$$

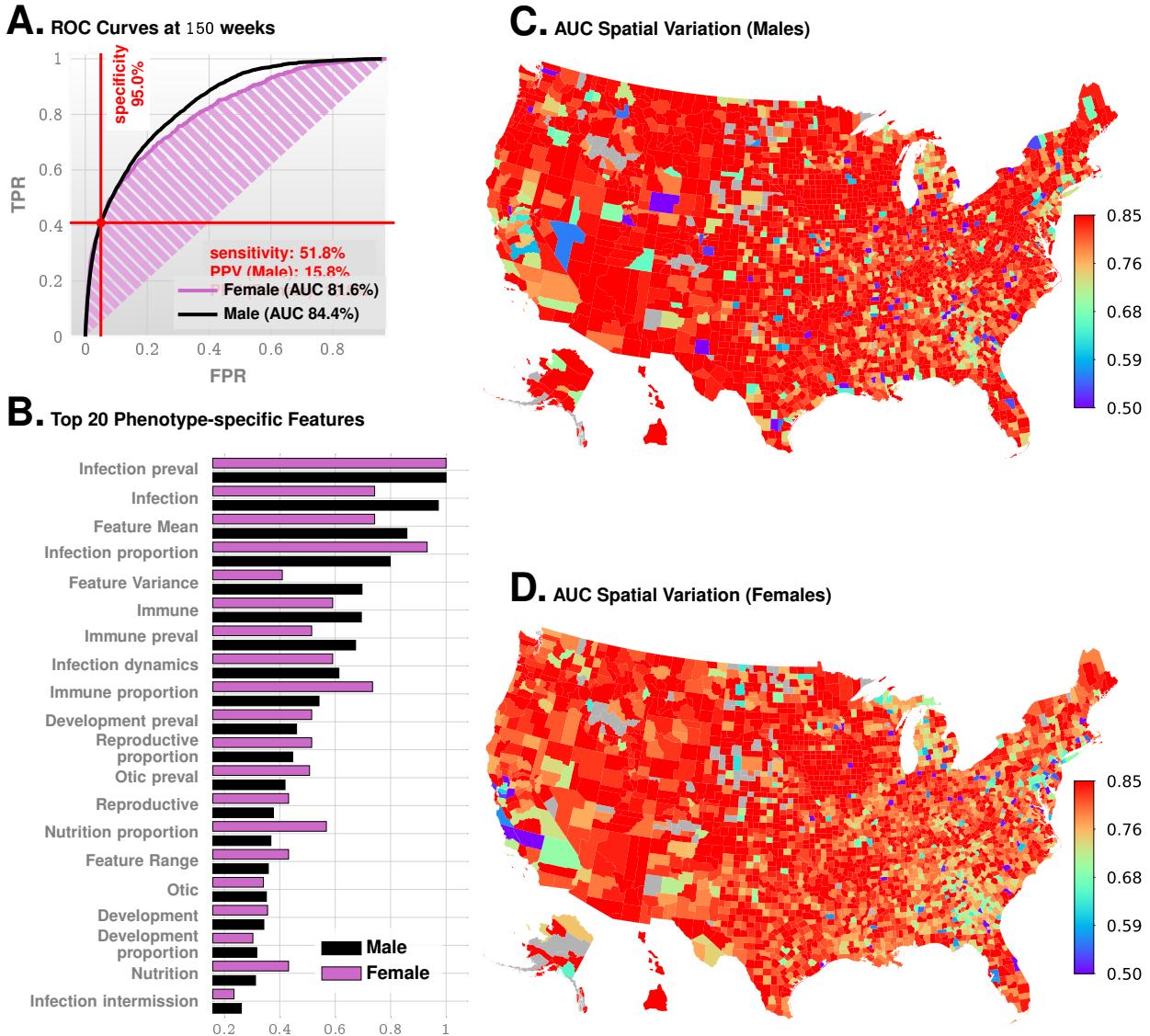
$$\rho = \frac{P}{N + P} \quad (12)$$



SI-Fig. 8. **Co-morbidity Patterns** for mental disorders, vaccinations and health-service encounters.

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**.



SI-Fig. 9. 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. 1 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 SI-Fig. 8A, the psychiatric codes all increase risk, and the vaccination codes (See SI-Fig. 8B) 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.

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:

$$\text{PPV} = \frac{t_p/P}{t_p/P + (f_p/N)(N/P)} = \frac{\text{TPR}}{\text{TPR} + ((N - t_n)/N)(N/P)} \quad (13)$$

$$\Rightarrow \text{PPV} = \frac{s}{s + (1 - c)(\frac{1}{P} - 1)} \quad (14)$$

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:

$$\text{PPV} = s/(s + 2.45) \leq 1/3.45 \sim 29\% \quad (15)$$

**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. 1C 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 study (14).

## VI. 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:

$$\text{TPR} = \frac{t_p}{P} = \frac{t'_p}{P'} \quad (16)$$

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

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

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. 6C, 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 (14). 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 VIII-A. First, we establish that our pipeline does not suffer from some important pitfalls arising in the work-flows associated with ASD diagnosis, and how the diagnostic codes in Electronic Health Records (EHR) are generated.

## VII. NOTE ON ASD CLINICAL DIAGNOSIS & UNCERTAINTY OF EHR RECORD

With no precise laboratory test for ASD, most families experience the following sequence of events (9, 13, 98): 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 (89, 101). 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 (9, 15).

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 function (9).

### A. 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 evaluation (98), 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) (9). 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) (82) and the Childhood Autism Rating Scale, Second Edition (CARS-2) (79). Current guidance from the American Academy of Pediatrics (9) 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 ASD (83). 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 judgment (83). 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” (83). 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 work-flows, 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. 6B illustrate that we have very little change in out-of-sample predictive performance, thus alleviating this concern.

#### *B. 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) (9). 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.

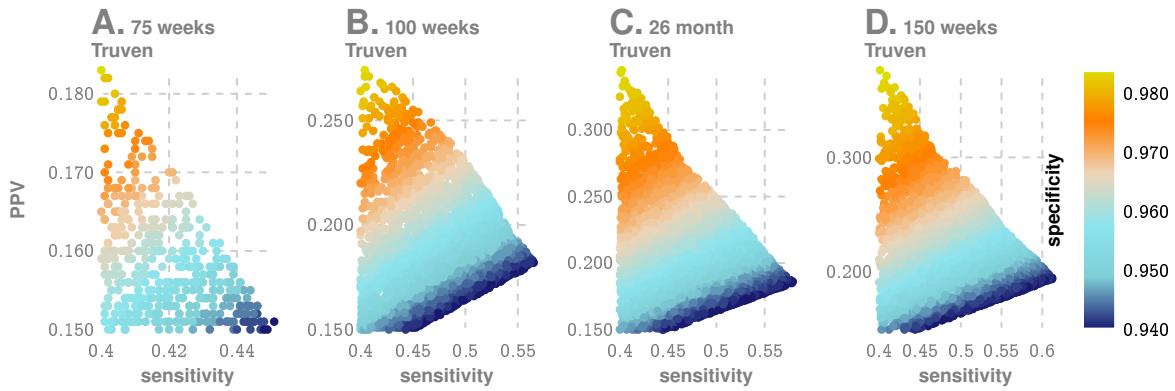
#### *C. Performance Comparison With M-CHAT/F*

The M-CHAT/F is the most studied and widely used tool for screening toddlers for ASD (9, 15).

Guthrie *et al.* (14) 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 IIa 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 screens (9). Additionally, while M-CHAT/F is relatively easy and quick to administer, the issue of time and resource commitment cannot be ignored (9). 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 AAP (14).



**SI-Fig. 10. 4D Search To Take Advantage of Data on Population Stratification (Using Prevalence of 2.23% as reported by CHOP (14)).** 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.

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.

### VIII. IMPROVING WAIT-TIMES FOR DIAGNOSTIC EVALUATIONS BY REDUCING FALSE POSITIVES

While children with ASD can be diagnosed as toddlers (91, 94) (developmental concerns may show up before the first birthday (73, 93)), the mean age of diagnosis is over 4 years (72). 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 evaluations (10), cost of care (11), lack of providers (12), lack of comfort in diagnosing by primary care providers (12), and other challenges, are all responsible to varying degrees that culminate in these delays (13), 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 (13, 37).

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.

#### A. 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  $\rho_i$  respectively, with  $i \in \{1, \dots, m\}$ . Let  $\beta_i$  be the relative size of the sub-populations. Thus, we have:

$$P = \sum_i P_i \quad (17)$$

$$N = \sum_i N_i \quad (18)$$

$$\beta_i = \frac{N_i + P_i}{N + P} \quad (19)$$

$$\rho_i = \frac{P_i}{N_i + P_i} = \frac{P_i}{\beta_i(N + P)} \quad (20)$$

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

$$s = t_p/P = \frac{\sum_i t_p|_i}{P} = \frac{\sum_i (t_P|_i/P_i) \times (\beta_i \rho_i (P + N))}{P} = \sum_i s_i \beta_i \frac{\rho_i}{\rho} \quad (21)$$

Thus, we end up with:

$$s = \sum_{i=1}^m s_i \gamma_i \quad (22a)$$

$$c = \sum_{i=1}^m c_i \gamma'_i \quad (22b)$$

$$PPV = \frac{s}{s + (1 - c)(\frac{1}{\rho} - 1)} \quad (22c)$$

where we have denoted:

$$\gamma_i = \beta_i \frac{\rho_i}{\rho}, \text{ and } \gamma'_i = \beta_i \frac{1 - \rho_i}{1 - \rho} \quad (22d)$$

Now, using Table III, we can compute the values for  $\gamma_i, \gamma'_i$ , as shown below.

SI-Table II  
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	$\geq 8$ 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 III  
POPULATION STRATIFICATION RESULTS ON LARGE M-CHAT/F STUDY(N=20,375) (14)

Id	Sub-population	Test Result	ASD positive	ASD Negative	Total %
A	M-CHAT/F $\geq 8$	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 IV  
 $\gamma, \gamma'$  COMPUTED FROM POPULATION STRATIFICATION RECORDED IN M-CHAT/F STUDY (14) ( $\rho = 0.0223$ )

Id	Sub-population	Test Result	$\beta_i$	$\rho_i$	$\gamma_i$	$\gamma'_i$
A	M-CHAT/F $\geq 8$	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

Using the prevalence and stratification parameters calculated from the CHOP study (See main text Table IV (14)), 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. (22). The results are shown in Fig. 10, 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. 1 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. 10 of this document, and in main text Fig. 3 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. (22) and hence the results are dependent on the population prevalence  $\rho$ . We report the dependence of the solution to the 4D optimization for population prevalence between 1.7% (CDC estimate (9)), and 2.23% (CHOP estimate (14)). In particular, it is illuminating to compare these results directly with M-CHAT/F performance, as shown in Fig. 3, 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%.

## IX. GENERATING PFSA MODELS FROM SET OF INPUT STREAMS WITH VARIABLE INPUT LENGTHS

Our PFSA reconstruction algorithm (31) 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.

## X. PROBABILISTIC FINITE STATE AUTOMATA INFERENCE

### A. Probabilistic Finite-State Automaton

Let  $\Sigma$  be a finite alphabet of symbols with size  $|\Sigma|$ . The set of sequences of length  $d$  over  $\Sigma$  is denoted by  $\Sigma^d$ . The set of finite but unbounded sequences over  $\Sigma$  is denoted by  $\Sigma^*$ , the Kleene star operation (87), i.e.  $\Sigma^* = \bigcup_{d=0}^{\infty} \Sigma^d$ . We use lower case Greek, for example  $\sigma$  or  $\tau$ , for symbols in  $\Sigma$ , and lower case Latin, for example  $x$  or  $y$ , for sequences of symbols, i.e.  $x = \sigma_1 \sigma_2 \dots \sigma_n$ . We use  $|x|$  to denote the length of  $x$ . The empty sequence is denoted by  $\lambda$ .

We denote the set of strictly infinite sequences over  $\Sigma$  by  $\Sigma^\omega$ , and the set of strictly infinite sequences having  $x$  as prefix by  $x\Sigma^\omega$ . Let  $\mathcal{S} = \{x\Sigma^\omega : x \in \Sigma^*\} \cup \{\emptyset\}$ , we can verify that  $\mathcal{S}$  is a semiring (92) over  $\Sigma^\omega$ . We use  $\mathcal{F}$  to denote the sigma algebra generated by  $\mathcal{S}$ .

**Definition 2** (Stochastic Process over  $\Sigma$ ). *A stochastic process over a finite alphabet  $\Sigma$  is a collection of  $\Sigma$ -valued random variables  $\{X_t\}_{t \in \mathbb{N}}$  indexed by positive integers (81).*

We are specifically interested in processes in which the  $X_t$ 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  $\mu_x$  induced by a sequence  $x \in \Sigma^*$  is the extension (92) to  $\mathcal{F}$  of the premeasure defined on the semiring  $S$  given by

$$\forall x, y \in \Sigma^*, \mu_x(y\Sigma^\omega) \triangleq \frac{Pr(xy)}{Pr(x)}, \text{ if } Pr(x) > 0 \quad (23)$$

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

**Definition 4** (Probabilistic Nerode Equivalence and Causal States (77)). 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 (78). 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,  $\Sigma$  is a finite alphabet,  $\delta : Q \times \Sigma \rightarrow \Sigma$  is called the transition map, and  $\tilde{\pi} : Q \rightarrow \mathbf{P}_\Sigma$ , where  $\mathbf{P}_\Sigma$  is the space of probability distributions over  $\Sigma$ , is called the transition probability. The entry of  $\tilde{\pi}(q)$  indexed by  $\sigma$  is denoted by  $\tilde{\pi}(q, \sigma)$ .

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

$$\pi_{q,q'} \triangleq \sum_{\{\sigma \in \Sigma | \delta(q, \sigma) = q'\}} \tilde{\pi}(q, \sigma) \quad (24)$$

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

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

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

$$\delta(q, x) \triangleq \delta(\delta(q, \sigma_1 \dots \sigma_{k-1}), \sigma_k) \quad (25)$$

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

$$\tilde{\pi}(q, x) \triangleq \prod_{i=1}^k \tilde{\pi}(\delta(q, \sigma_1 \dots \sigma_{i-1}), \sigma_i) \quad (26)$$

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 connected (74). 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 (90, 100) of  $G$  and is denoted as  $\wp_G$ , or  $\wp$  if  $G$  is understood.

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

$$\Gamma_\sigma|_{q,q'} \triangleq \begin{cases} \tilde{\pi}(q, \sigma) & \text{if } \delta(q, \sigma) = q', \\ 0 & \text{if otherwise.} \end{cases} \quad (27)$$

$\Gamma_\sigma$  is called event-specific transition matrix, with the event being that  $\sigma$  is current the output.  $\Gamma_\sigma$  can also be extended to arbitrary  $x \in \Sigma^*$  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  $\wp_0$  on  $Q$ , the distribution on  $Q$  induced by a sequence  $x$  is given by  $\wp_{G, \wp_0}^T(x) = [\wp_0^T \Gamma_x]$  with  $\wp_{G, \wp_0}(\lambda) = \wp_0$ . The entry indexed by  $q \in Q$  of the vector  $\wp_{G, \wp_0}(x)$  is written as  $\wp_{G, \wp_0}(x, q)$ . When  $\wp_0 = \wp_G$ , the stationary distribution of  $G$ , we write  $\wp_{G, \wp_0}(x)$  as  $\wp_G(x)$ , or simply as  $\wp(x)$ , if  $G$  is understood.

**Definition 11** (Stochastic Process Generated by a PFSA). Let  $G = (Q, \Sigma, \delta, \tilde{\pi})$  be a PFSA and let  $\wp_0$  be a distribution on  $Q$ , the  $\Sigma$ -valued stochastic process  $\{X_t\}_{t \in \Sigma}$  generated by  $G$  and  $\wp_0$  satisfies that  $X_1$  follows the distribution  $\wp_0$  and  $X_{t+1}$  follows the distribution  $\wp_{G, \wp_0}(X_1 \dots X_t)$  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  $\rho_G$ , 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  $\phi_\lambda$  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 state (99). 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  $\varepsilon$ -synchronizing sequence to a state  $q \in Q$  if

$$\|\rho_G(x) - e_q\|_\infty \leq \varepsilon. \quad (28)$$

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

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

$$\Sigma_{\varepsilon, (q, r)}^d \triangleq \{x \in \Sigma^d : x \text{ jointly } \varepsilon\text{-synchronizing to } (q, r)\} \quad (29)$$

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

$$p_G(q, r) \triangleq \lim_{d \rightarrow \infty} p_G\left(\Sigma_{\varepsilon, (q, r)}^d\right) \quad (30)$$

*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*

$$Q_c \triangleq \{(q, r) \in Q_G \times Q_H : (q, r) \text{ is a } G\text{-joint pair}\} \quad (31)$$

The inference algorithm for PFSA is called **GenESeSS** for Generator Extraction Using Self-similar Semantics. With an input sequence  $x$  and a hyperparameter  $\varepsilon$ , **GenESeSS** outputs a PFSA in the following three steps: 1) approximate an almost synchronizing sequence; 2) identify the transition structure of the PFSA; 3) calculate the transition probabilities of the PFSA. See Alg. 1 for detail.

## XI. 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 PFSA(80). Similarly, by KL divergence of two PFSA, we mean the KL divergence between the two processes generated by them (96). More precisely, we have*

$$\mathcal{H}(G) = - \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{x \in \Sigma^d} p(x) \log p(x) \quad (32)$$

*and the KL divergence*

$$\mathcal{D}_{KL}(G \parallel H) = \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{x \in \Sigma^d} p_G(x) \log \frac{p_G(x)}{p_H(x)} \quad (33)$$

*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*

$$\mathcal{H}(G) = \sum_{q \in Q} \rho_G(q) \cdot h(\tilde{\pi}(q)) \quad (34)$$

*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  $\mu_G$  being absolutely continuous with*

**Algorithm 1: GenESeSS**


---

**Data:** A sequence  $x$  over alphabet  $\Sigma$ ,  $0 < \varepsilon < 1$   
**Result:** State set  $Q$ , transition map  $\delta$ , and transition probability  $\tilde{\pi}$

```

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

- 1 Let  $L = \lceil \log_{|\Sigma|} 1/\varepsilon \rceil$ ;
- 2 Calculate the **derivative heap**  $D_\varepsilon^x$  equaling  $\{ \hat{\phi}_y^x : y \text{ is a sub-sequence of } x \text{ with } |y| \leq L \}$ ;
- 3 Let  $\mathcal{C}$  be the convex hull of  $D_\varepsilon^x$ ;
- 4 Select  $x_0$  with  $\hat{\phi}_{x_0}^x$  being a vertex of  $\mathcal{C}$  and has the highest frequency in  $x$ ;

```

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

- 5 Initialize  $Q = \{q_0\}$ ;
- 6 Associate to  $q_0$  the **sequence identifier**  $x_{q_0}^{\text{id}} = x_0$  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  $q_{\text{new}}$  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  $\delta$ , and denote the vertex set of the subgraph again by  $Q$ ;

```

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

- 20 Initialize counter  $N[q, \sigma]$  for each pair  $(q, \sigma) \in Q \times \Sigma$ ;
- 21 Choose a random starting state  $q \in Q$ ;
- 22 **for**  $\sigma \in x$  **do**
- 23 Let  $N[q, \sigma] = N[q, \sigma] + 1$ ;
- 24 Let  $q = \delta(q, \sigma)$ ;
- 25 Let  $\tilde{\pi}(q) = \llbracket (N[q, \sigma])_{\sigma \in \Sigma} \rrbracket$ ;
- 26 **return**  $Q, \delta, \tilde{\pi}$ ;

---

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

$$\mathcal{D}_{KL}(G \parallel H) = \sum_{(q, r) \in Q_c} p_G(q, r) D_{KL}(\tilde{\pi}_G(q) \parallel \tilde{\pi}_H(r)) \quad (35)$$

**Definition 18** (Log-likelihood). Let  $x \in \Sigma^d$ , the log-likelihood (80) of a PFSA  $G$  generating  $x$  is given by

$$L(x, G) = -\frac{1}{d} \log p_G(x) \quad (36)$$

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

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

$$L(x, H) \rightarrow \mathcal{H}(G) + \mathcal{D}_{KL}(G \parallel H) \quad (37)$$

in probability as  $d \rightarrow \infty$ .

*Proof.* We first notice that

$$\sum_{x \in \Sigma^d} p_G(x) \log \frac{p_G(x)}{p_H(x)} = \sum_{x \in \Sigma^{d-1}} \sum_{\sigma \in \Sigma} p_G(x) \wp_G(x) \left. \tilde{\Pi}_G \right|_\sigma \log \frac{p_G(x) \wp_G(x) \left. \tilde{\Pi}_G \right|_\sigma}{p_H(x) \wp_H(x) \left. \tilde{\Pi}_H \right|_\sigma} \quad (38)$$

$$= \sum_{x \in \Sigma^{d-1}} p_G(x) \log \frac{p_G(x)}{p_H(x)} + \underbrace{\sum_{x \in \Sigma^{d-1}} p_G(x) \sum_{\sigma \in \Sigma} \wp_G(x) \left. \tilde{\Pi}_G \right|_\sigma \log \frac{\wp_G(x) \left. \tilde{\Pi}_G \right|_\sigma}{\wp_H(x) \left. \tilde{\Pi}_H \right|_\sigma}}_{D_d} \quad (39)$$

**Algorithm 2:** Log-likelihood

---

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

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

---

By induction, we have  $\mathcal{D}_{\text{KL}}(G \parallel H) = \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=1}^d D_i$ , and hence by Cesàro summation theorem (85), we have  $\mathcal{D}_{\text{KL}}(G \parallel H) = \lim_{d \rightarrow \infty} D_d$ . 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

$$-\frac{1}{n} \sum_{i=1}^n \log \varphi_H(x^{[i-1]}) \tilde{\Pi}_H \Big|_{\sigma_i} = \underbrace{\frac{1}{n} \sum_{i=1}^n \log \frac{\varphi_G(x^{[i-1]}) \tilde{\Pi}_G \Big|_{\sigma_i}}{\varphi_H(x^{[i-1]}) \tilde{\Pi}_H \Big|_{\sigma_i}}}_{A_{x,n}} - \underbrace{\frac{1}{n} \sum_{i=1}^n \log \varphi_G(x^{[i-1]}) \tilde{\Pi}_G \Big|_{\sigma_i}}_{B_{x,n}} \quad (40)$$

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

$$\lim_{n \rightarrow \infty} A_{x,n} = \lim_{d \rightarrow \infty} D_d = \mathcal{D}_{\text{KL}}(G \parallel H) \quad (41)$$

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

## XII. PIPELINE OPTIMIZATION

### A. 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` algorithm (31) (See Section X), 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.

### B. Algorithms

The key data processing approach is outlined in Algorithm 3. The remaining steps of the approach are sketched in Algorithm 4. Fig. 11 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.

## XIII. EXAMPLE RUN WITH RELEASED APPLICATION

### A. Prerequisites & Installation

The minimum prerequisites for running ehrzero are the following:

- 1) A x64 system running any flavor of Linux.

**Algorithm 3:** 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;

```

---

- 2) A working python 3.x installation  
 3) scikit-learn, version = 0.20.0

Installation:

```
pip3 install ehrzero --user
```

#### B. EHR data format

Diagnostic data stored in text file, one line per patient as follows: patient id, gender, and list of space-separated, comma-delineated 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
```

#### C. 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.

#### D. Sample Python script risk estimation

The script version is similar to the one mentioned before.

**Algorithm 4:** Prediction Pipeline Training

---

```

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

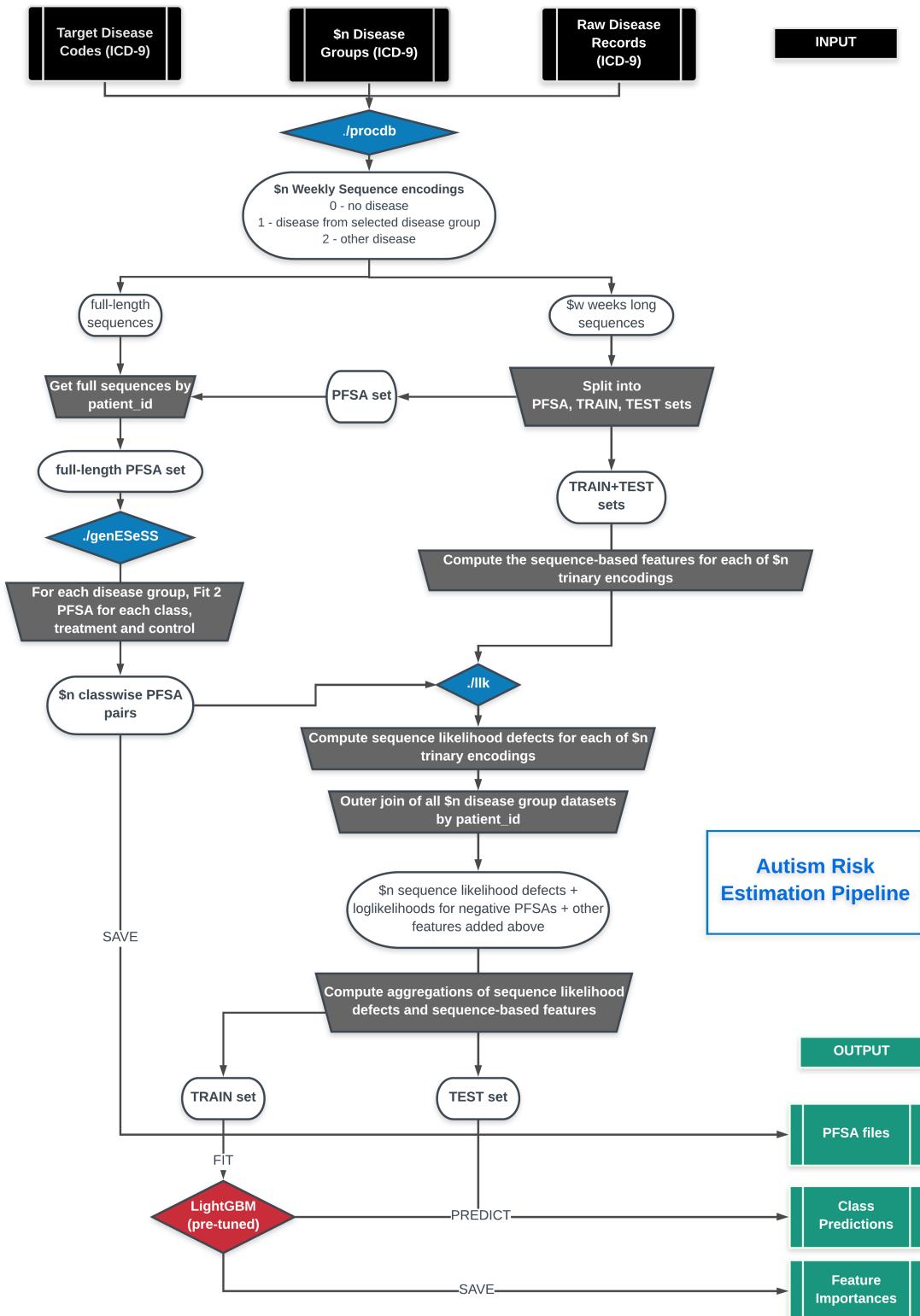
1 DiseaseDatasets ← new Dictionary();
2 for Dataset, DiseaseGroup ∈ zip(Encoding, DiseaseGroups) do
3   PFSAsset, LLset ← TrainTestSplit(Dataset, w.r.t = "target");
4   df ← new Dataframe();
5   df[patientID] ← LLset[patientID];
6   df[target] ← LLset[target];
7   //Generate 2 PFSA for each class;
8   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;

```

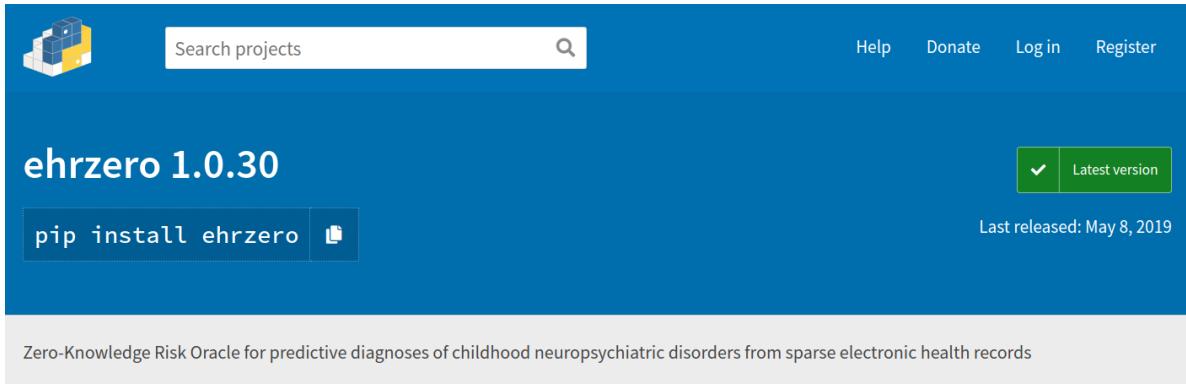
---

Once ehrzero package is installed, locate its directory and go to `textt..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
```



SI-Fig. 11. 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 key algorithms here are genESeSS (31) and the llk which does the sequence likelihood computation described in Section XI



Navigation

## Project description

Project description

downloads 174/month

pypi v1.0.30

Release history

SI-Fig. 12. Screen capture of the page on pypi.org hosting the released application Link: <https://pypi.org/project/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	AA Abby	200	0.000174	0.028200	0.920977
0	AA Abby	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. 13. Python code prediction example

```
(base) [onishchenko@midway2-login1 example]$ cat tests/ZERO_example.dat
M:44 380.10:101 381.81:11 084.6
M:9 380.10:104 381.81:11 084.6
M:99 380.10:104 381.81:11 084.6:98 380.11
M:9 380.10:104 381.71:11 084.6
M:9 390.11:4 391.71:11 084.6
(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
A000000001 100 0.001703 0.098002 68.54
A000000001 200 0.001834 0.105517 72.43
A000000001 300 0.001834 0.105517 72.43
A000000002 100 0.002030 0.116816 60.94
A000000002 200 0.001426 0.082033 80.01
A000000002 300 0.001426 0.082033 80.01
A000000003 100 0.001703 0.098002 68.54
A000000003 200 0.001706 0.098170 74.53
A000000003 300 0.001853 0.106625 72.18
A000000004 100 0.000512 0.029449 94.86
A000000004 200 0.001450 0.083436 79.53
A000000004 300 0.001450 0.083436 79.53
A000000005 100 0.000658 0.037868 91.73
A000000005 200 0.000658 0.037868 95.21
A000000005 300 0.000658 0.037868 95.21
(base) [onishchenko@midway2-login1 example]$
```

SI-Fig. 14. Python script prediction example