

# Computational Phenotyping in Electronic Health Records

*Slides from Pedro L. Teixeira, PhD*

# Evaluation of Data Sources and Algorithms for Hypertension Phenotyping in the Electronic Health Record

# Hypertension has high prevalence and mortality

- Consistently high blood pressure readings:
  - Systolic > 139 mmHg OR Diastolic > 89 mmHg
- Hypertension is a risk factor for:
  - Heart attack
  - Heart failure
  - Stroke
  - Kidney disease
  - Aneurysm
- Affects approximately 1/3<sup>rd</sup> of Americans
- Contributes to 1/6<sup>th</sup> of deaths
- Can be well controlled with lifestyle changes and medications

A seemingly simple question...

Which of the patients in the  
EHR have hypertension?

Phenotyping algorithms are usually sets of conditions applied across data types

- If the individual has three hypertensive blood pressure readings

**OR**

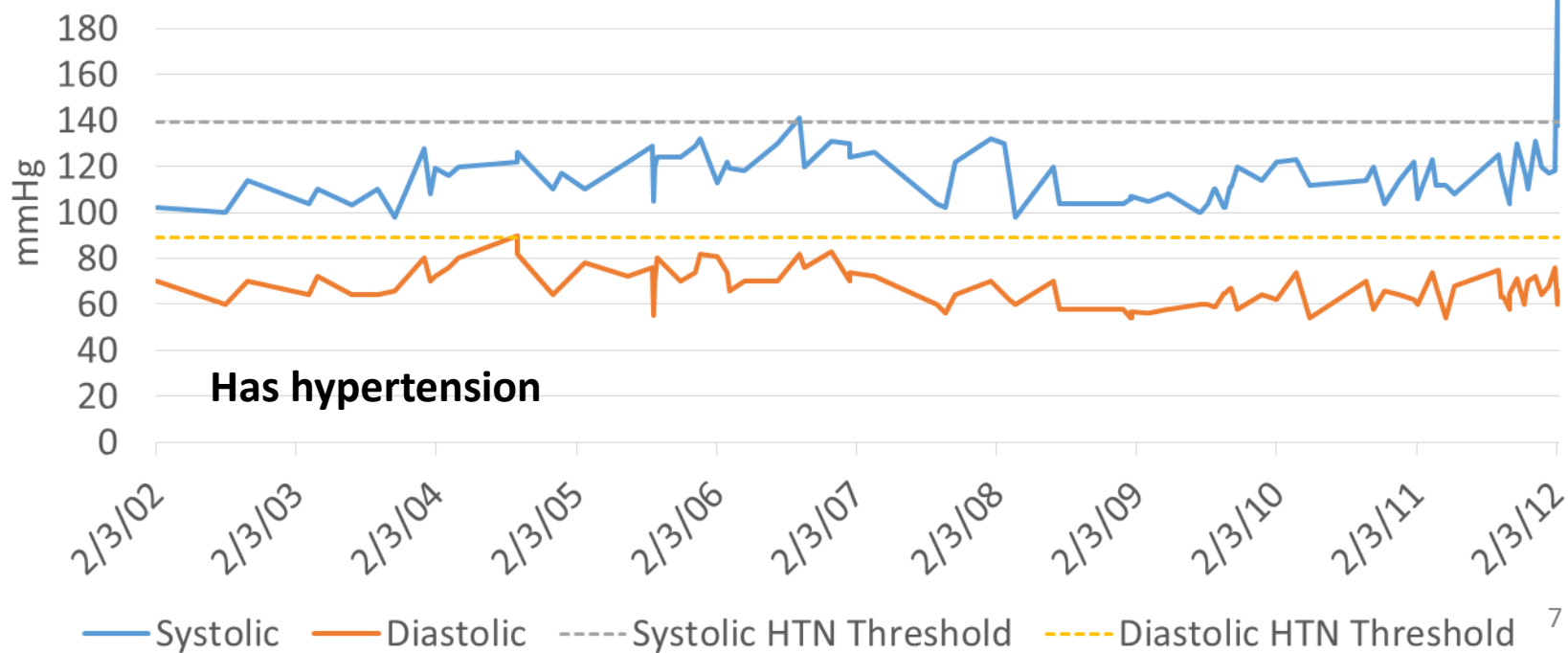
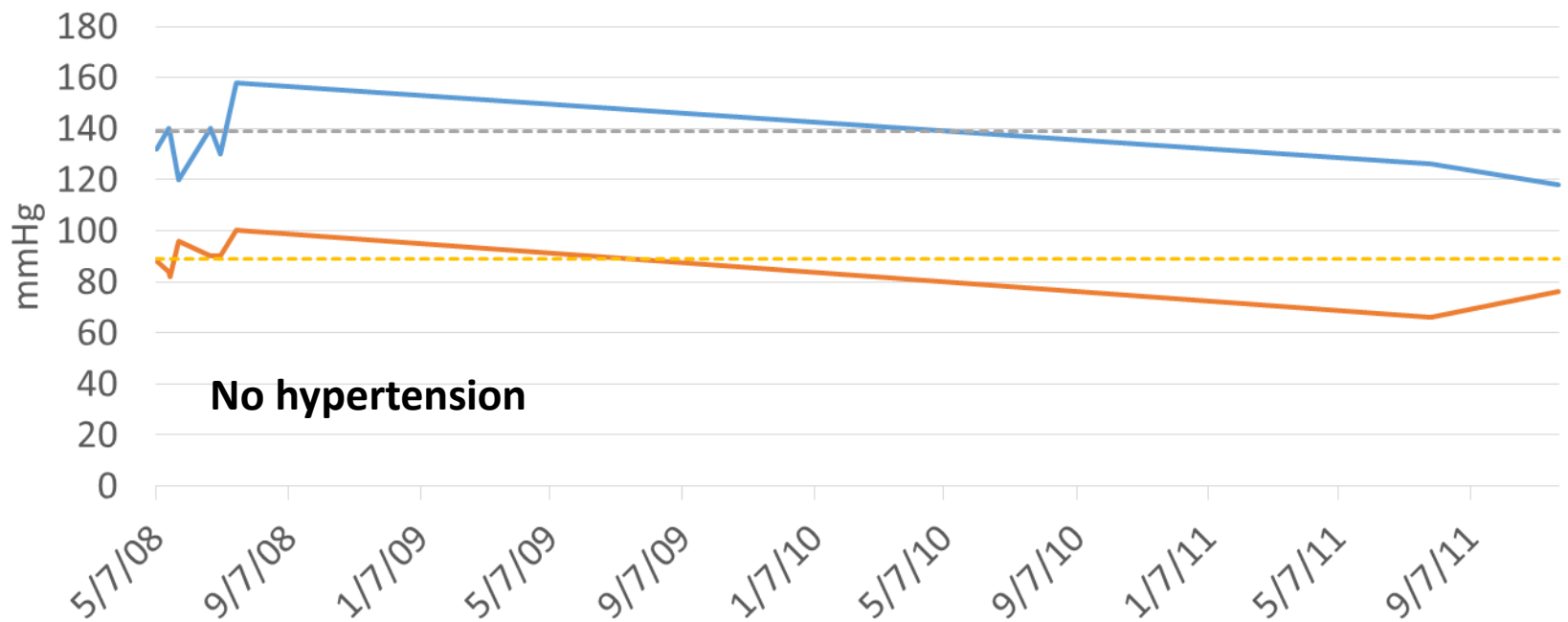
- (If the individual has a billing code for hypertension

**AND**

- One or more medications that are used to treat hypertension)

# Phenotyping algorithms are important for research and clinical practice

- Large scale genotype-phenotype studies
- Outbreak surveillance
- Clinical decision support
- Quality improvement



[illegible]



# Regular expressions can match text with limited sets of rules

- Search for 'hypertension' OR 'htn' – ignore case
- Matches:
  - Hypertension
  - hypertension
  - HTN
  - HtN
  - htn
  - ... pulmonary-hypertension
  - The patient does not have hypertension.
  - Mr. Smith's father and mother have hypertension.

# Advanced natural language processing tools can extract concepts from narrative text

Pre-processing	<p>Mr. Smith is a 57yo man with HLD and HTN. He's here today for f/u.</p>					
Sentence boundary detection	Mr. Smith is a 57yo man with HLD and HTN.   He's here today for f/u.					
Tokenization	Mr . Smith is a 57yo man with HLD and HTN . He's here today for f / u .					
Part-of-speech tagging	NNP NNP VBZ DT NN NN IN NNP CC NNP .   PRP ' VBZ RB NN IN NN .					
Named entity recognition	Mr . Smith	man	HLD	HTN	f / u	
	Name title	Male gender	Hyperlipidemia	Hypertension, Malignant	Follow-up	
	C01704446	C00024554	C00020473	[Disease/Finding]	C01522577	
				C00020540		

Denny, J., *et al* (2003). "Understanding" medical

school curriculum content using KnowledgeMap

Friedman C, Rindflesch TC, Corn M. Natural language processing: state of the art

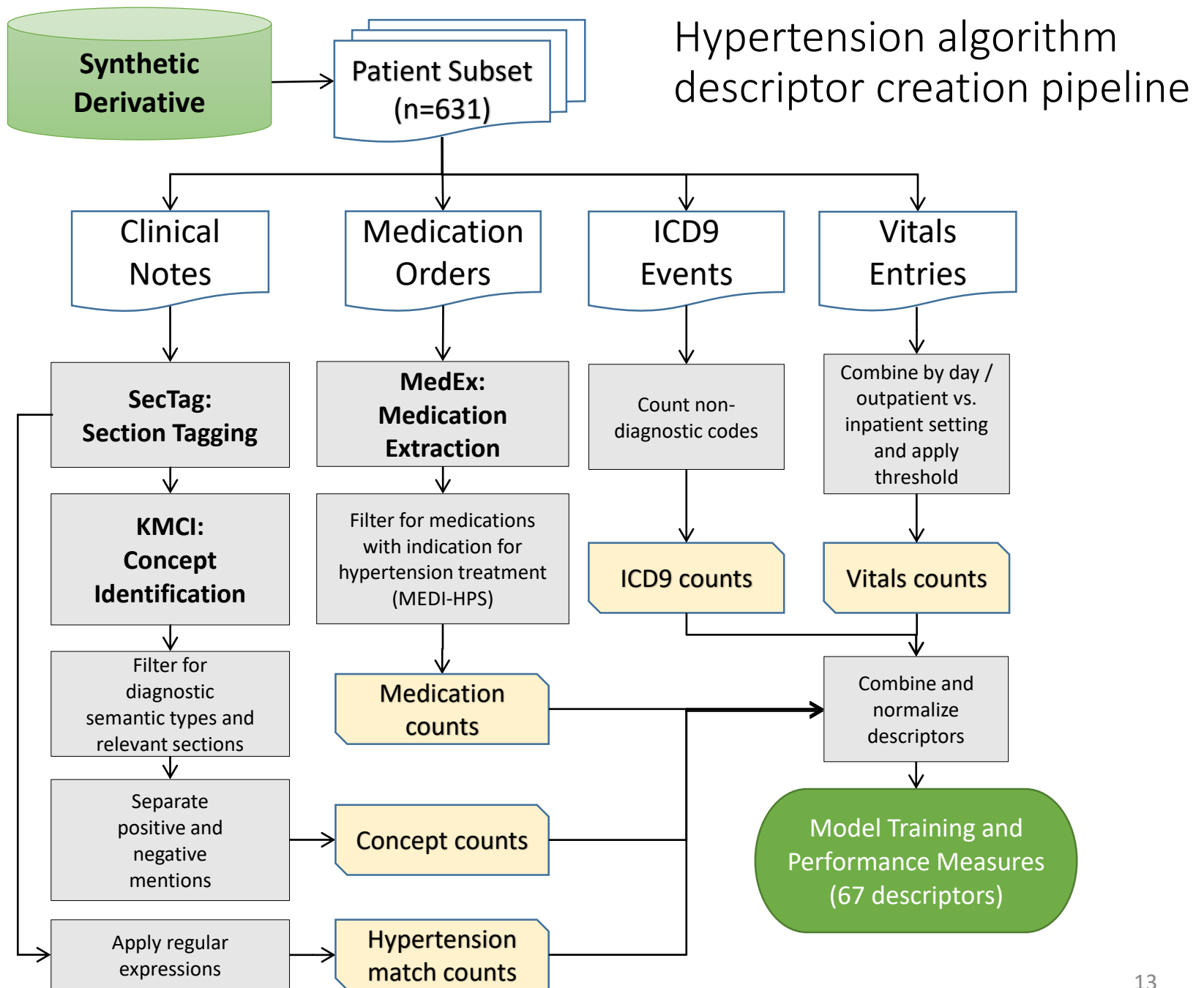
and prospects for significant progress, a workshop sponsored by the National Library of Medicine. 2013

# Goal: Evaluate various EHR data sources for hypertension phenotyping

- Examine four categories of data:
  - Vitals
  - Billing (ICD9) codes
  - Medication orders
  - Clinical notes
- Develop and test various algorithms on these data sources
  - Individually
  - Categorically
  - In combination
  - Using machine learning

# Hypertension Algorithm Development Dataset (n=643)

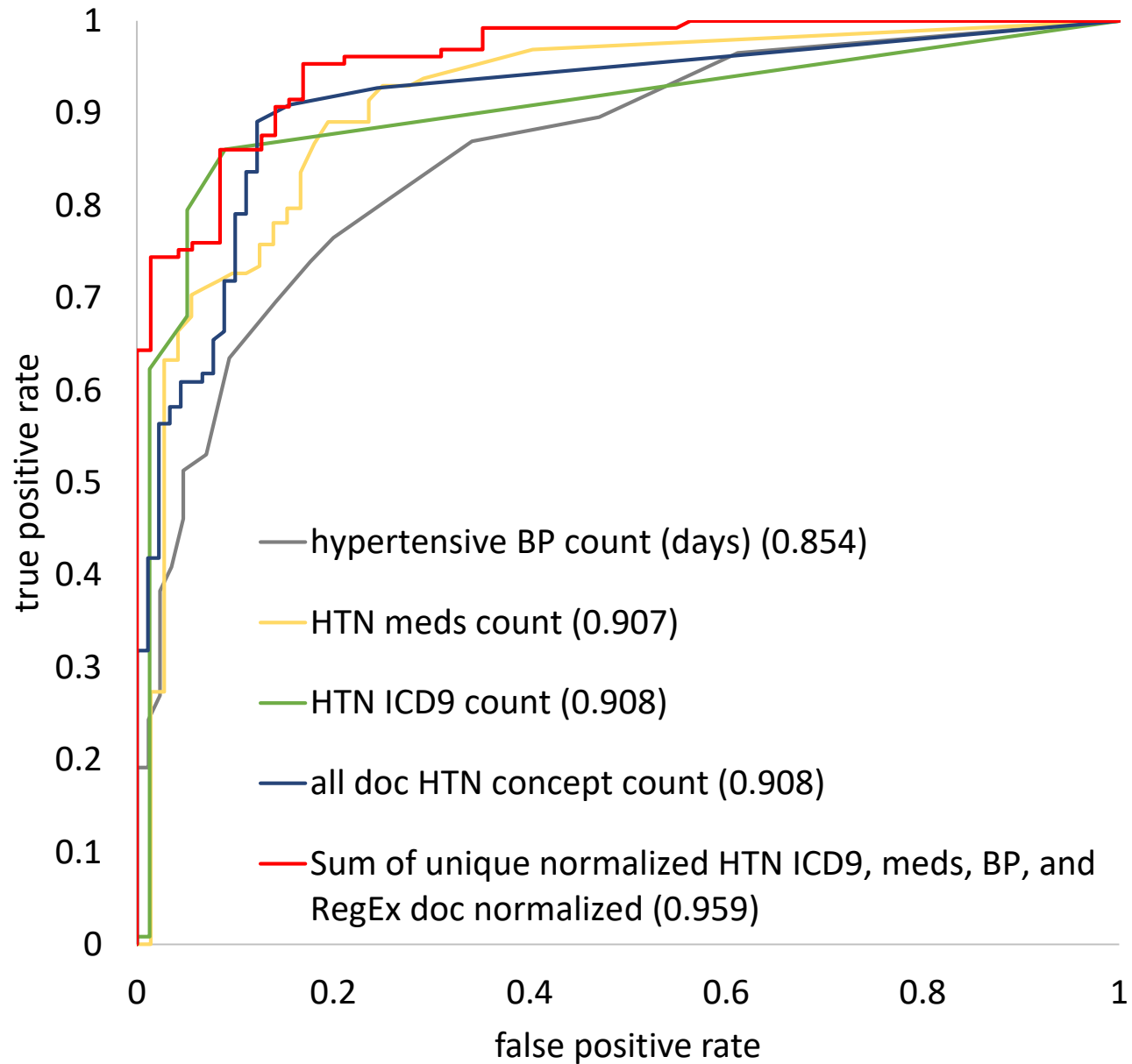
- 643 Randomly selected adults with:
  - 2 or more outpatient visits between 1/1/07 -1/1/09
  - 2 or more vitals readings between 1/1/07 -1/1/09
- 631 labeled as case or control
- Manually reviewed by three physicians and one medical student
  - 20% overlap (on initial 303 of 643)
  - kappa = 0.93
- $\sim 2/3^{\text{rd}}$  Cases



# All 67 descriptors extracted

Category	Column description	Category	Column description	Category	Column description
ICD9, medications, all BP	date-ICD9 count	All vitals with separate outpatient and inpatient blood pressures	outpatient visits with vitals count (days)	NLP-based concepts and normalized descriptors	PL concept count
	unique ICD9 count		outpatient hypertensive BP count (days)		DS concept count
	HTN ICD9 count		median systolic (outpatient)		HPC concept count
	unique HTN ICD9 count		median diastolic (outpatient)		All doc concept count
	HTN ICD9 count normalized		outpatient hypertensive BP count normalized		PL HTN concept count
	HTN ICD9 count unique normalized		outpatient vitals density		DS HTN concept count
	unique HTN ICD9 count normalized		inpatient visits with vitals count (days)		HPC HTN concept count
	unique HTN ICD9 count unique normalized		inpatient hypertensive BP count (days)		All doc HTN concept count
	meds count		median systolic (inpatient)		PL HTN concept count doc normalized
	unique meds count		median diastolic (inpatient)		DS HTN concept count doc normalized
	HTN meds count		inpatient hypertensive BP count normalized		HPC HTN concept count doc normalized
	unique HTN meds count		inpatient vitals density		All doc HTN concept count doc normalized
	HTN meds count normalized		median pulse (all)		PL HTN concept count concept normalized
	HTN meds count unique normalized		median pulse (outpatient)		DS HTN concept count concept normalized
	unique HTN meds count normalized		median pulse (inpatient)		HPC HTN concept count concept normalized
	unique HTN meds count unique normalized		PL count		All doc HTN concept count concept normalized
	max age	Document counts	DS count	Regular expression counts and normalized descriptors	
	vital reading time span (days)		HPC count		
	visits with vitals count (days)	Regular expression counts and normalized descriptors	All doc count		
	hypertensive BP count (days)		PL HTN RegEx doc matches		
	median systolic (all)		DS HTN RegEx doc matches		
	median diastolic (all)		HPC HTN RegEx doc matches		
	hypertensive BP count normalized		All doc HTN RegEx doc matches		
	vitals density		PL HTN RegEx doc matches doc normalized		
			DS HTN RegEx doc matches doc normalized		
			HPC HTN RegEx doc matches doc normalized		
			All doc HTN RegEx doc matches doc normalized		

# Combining categories improves performance



# Random Forests are a robust and easy-to-use machine learning method

## **Pros**

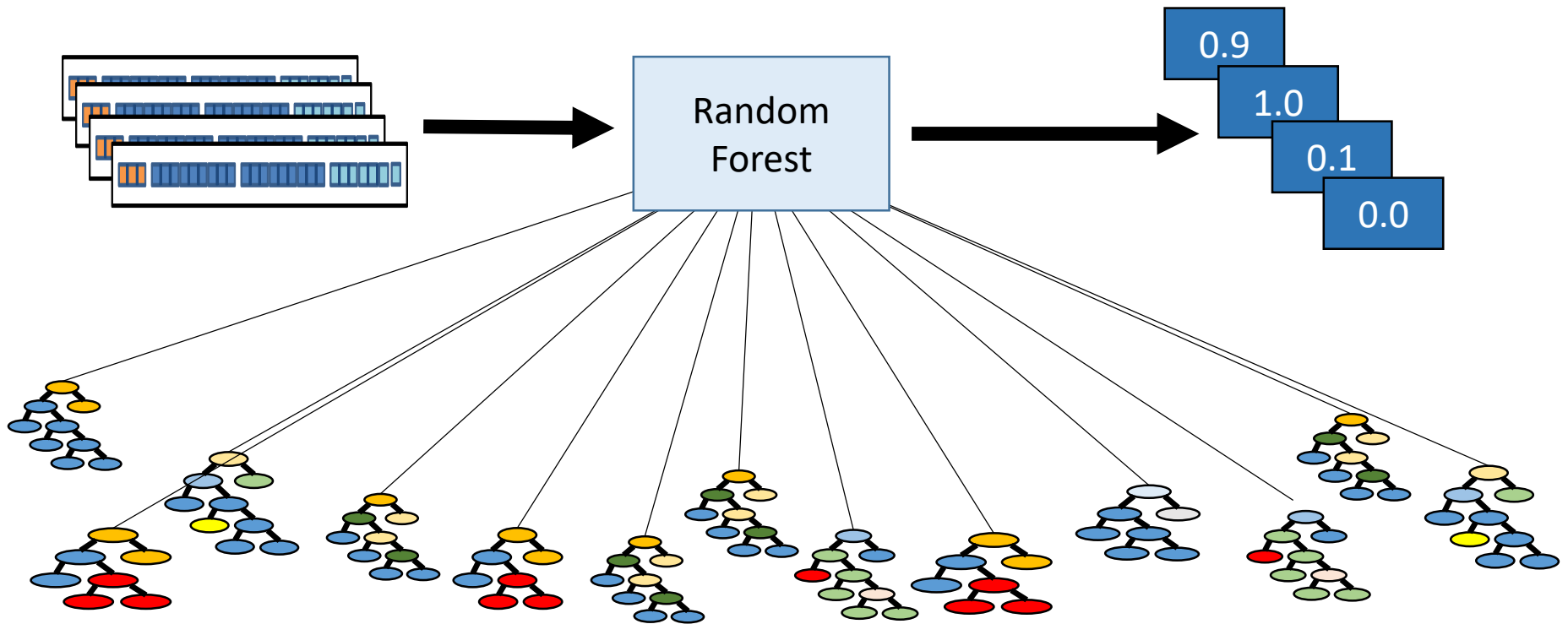
- More stable than decision trees
- Good on “messy” data
- Relatively few parameters to optimize
- Suggests descriptor “importance”

## **Cons**

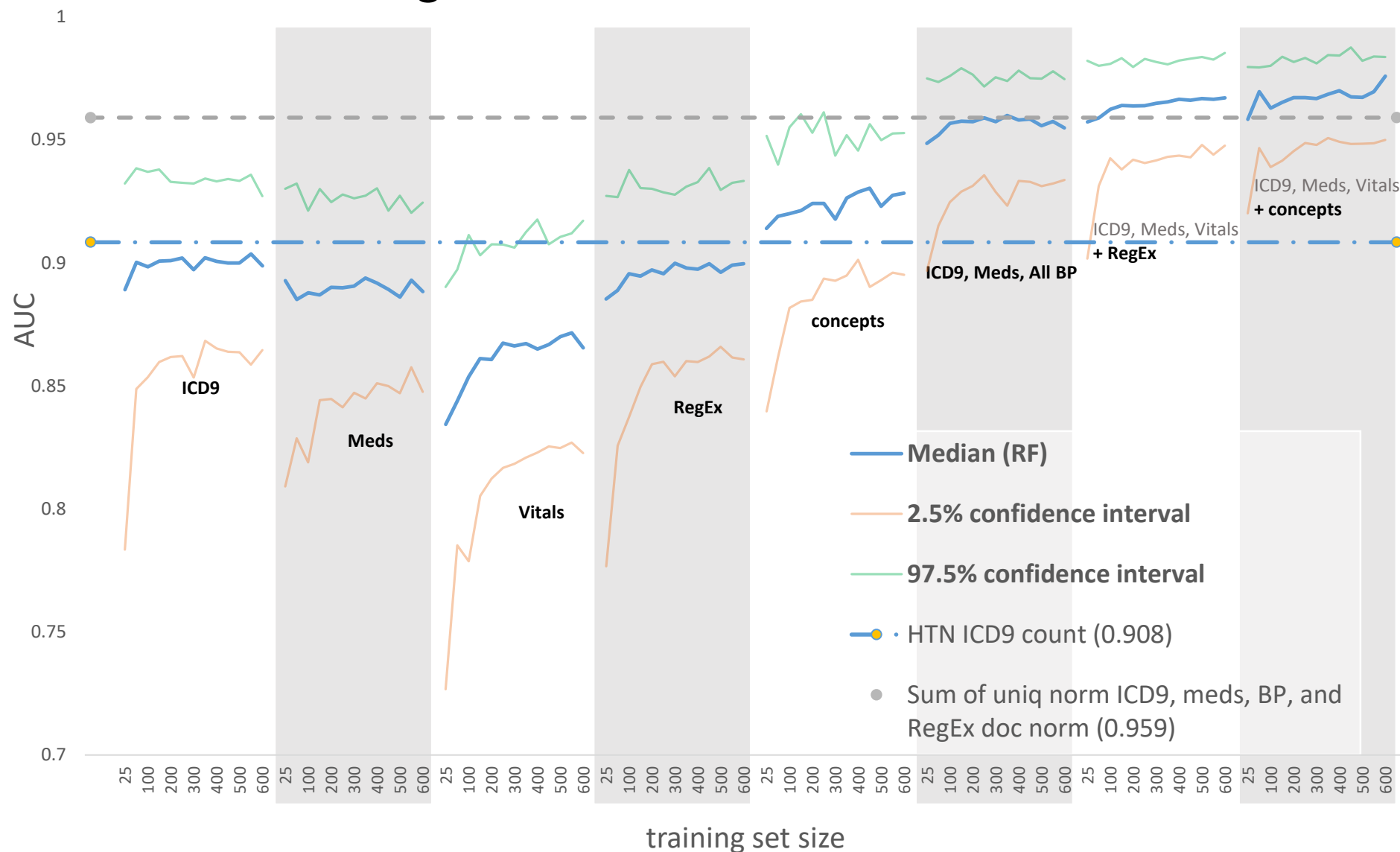
- More time intensive to train
- “Black box”
- May not capture complex relationships
- Random method



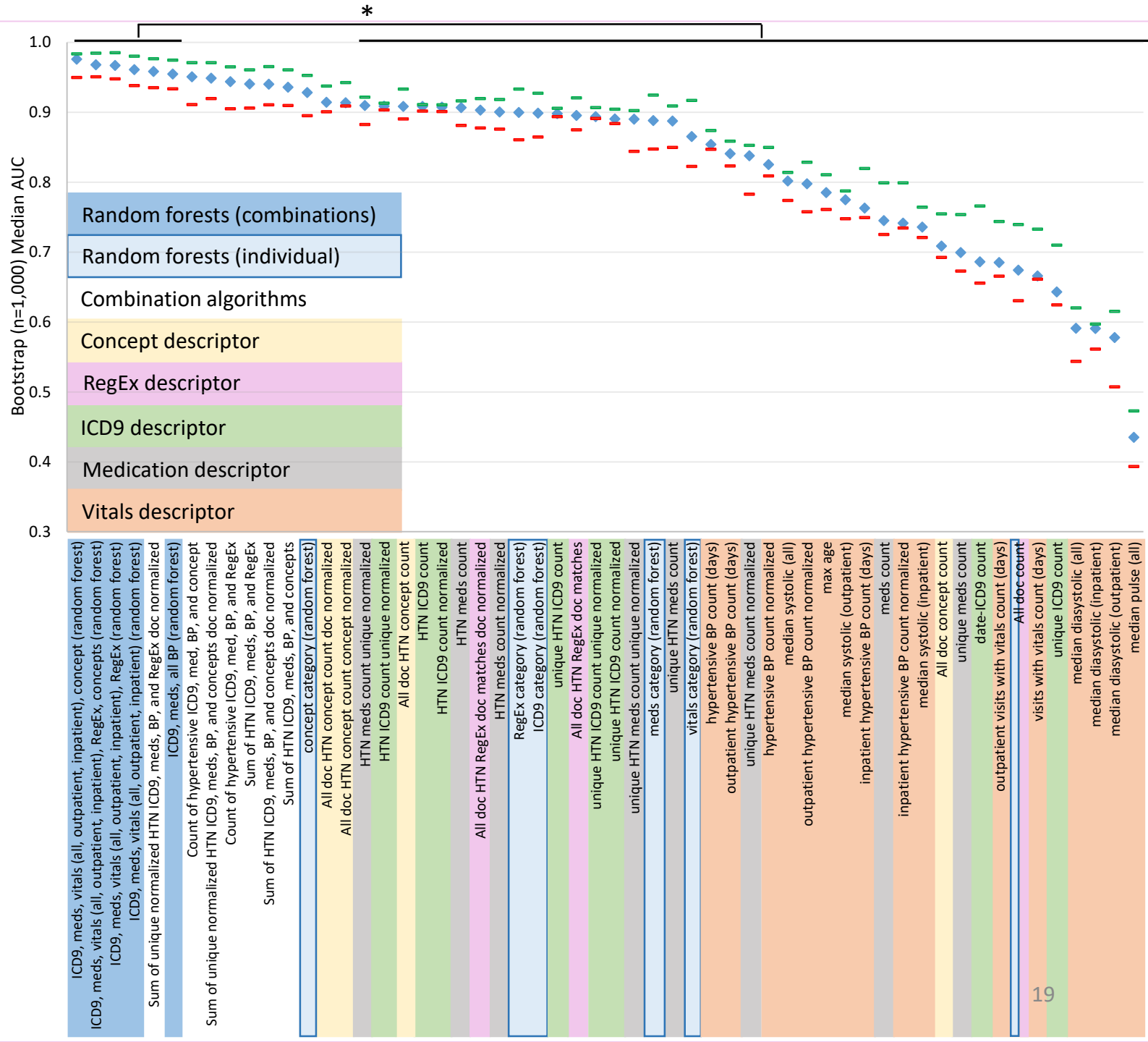
# Random forest result is the mode of the trees



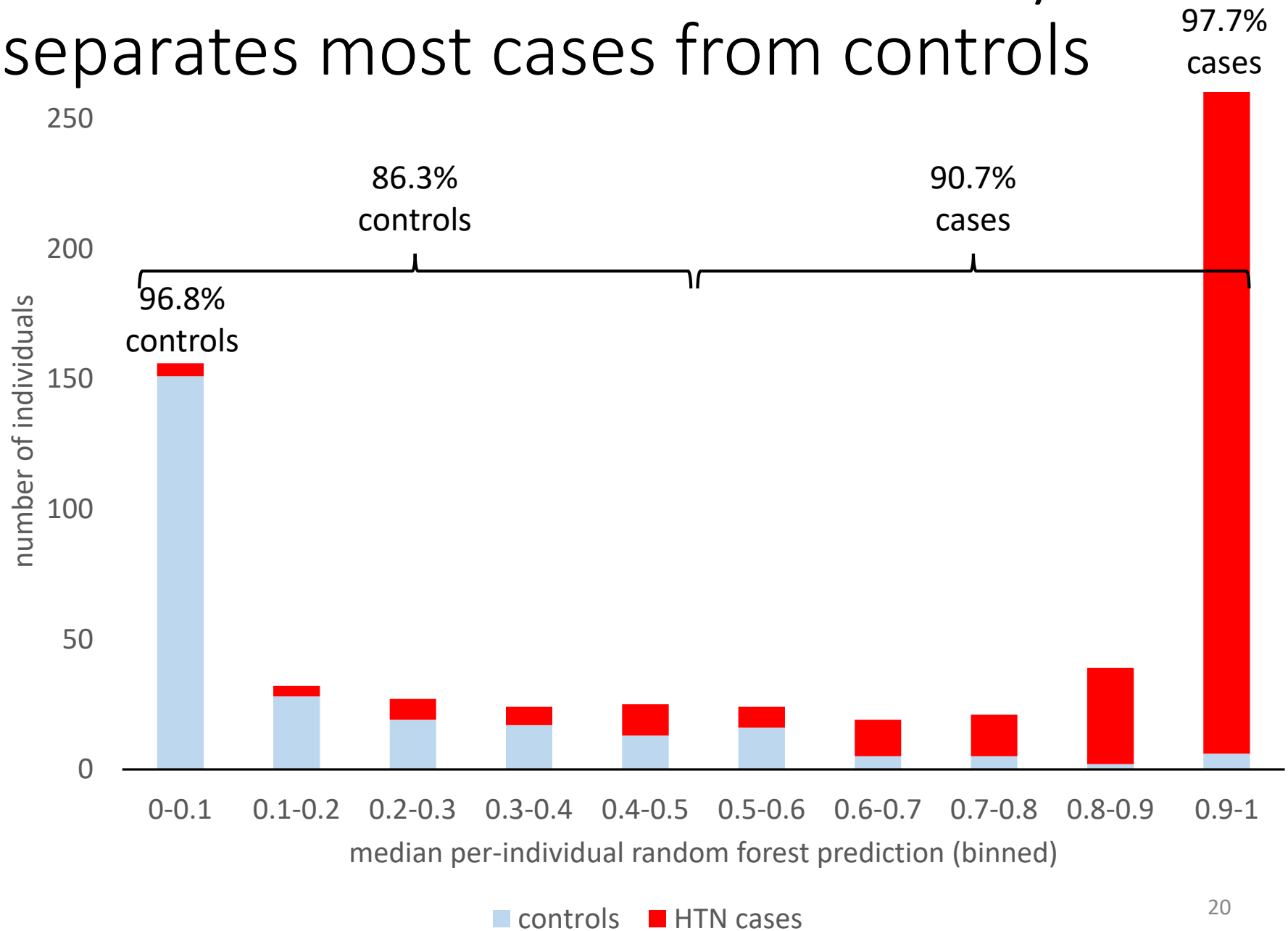
# Random forests trained on combinations of data achieve the highest AUC



*Vitals – Includes all blood pressure readings (inpatient, outpatient, and combined) and pulse  
BP - Includes all blood pressure readings (combined only)*



# The best random forest effectively separates most cases from controls



# We created a portable Konstanz Information Miner (KNIME) module

All nodes with yellow backgrounds require input (filepath locations)

Input the demographics (DOB) filepath

CSV Reader

Read from input file demographics (DOB)

Input the labels filepath (necessary if you'd like sensitivity, specificity, and PPV estimates for each model)

CSV Reader  
Labeled IDs

Specify the path for all output files

Table Creator

Specify path for outputs

Table Creator

Output filenames

Table Creator

Specify path to model

Table Creator

Input R library paths

Column Appender Java Snippet

Node 248

Add paths to each filename

Joiner

Node 250

String to Date/Time

Convert date strings to dates

Joiner

Combine IDs, DOB, and Labels (if available)

Specify the path to the provided F randomForest models (.rda files):  
Make sure to include final '/'

Specify the R library paths here (can be found from .libPaths() in R):

Also, make sure to install these packages  
install.packages("rJava")  
install.packages("randomForest")  
install.packages("ROCR")

Specify the filenames for each of your input files for (relevant nodes highlighted with yellow backgrounds):

- Demographics (DOB)
- Billing codes
- Medications
- Blood pressure readings
- Outpatient dates
- Pulse
- Document type counts
- Regular expression hypertension counts
- Natural language processing hypertension concept counts

CSV Reader

Read from input file Ensure normalization of multiple entries per day to one

GroupBy

All ICD9-day counts

GroupBy

All unique ICD9 counts

Table Creator

Hypertension ICD9 codes

Joiner

Inner join input ICD9 codes with HTN ICD9 Codes

GroupBy

Calculate filtered HTN ICD9 code counts

GroupBy

Calculate filtered HTN ICD9 unique code counts

Column Rename

Rename

ALL\_DATE\_ICD9\_COUNT

Joiner

Column Rename

Rename

ALL\_UNIQ\_ICD9\_COUNT

Joiner

Column Rename

Rename

FILT\_HTN\_ICD9\_COUNT

Column Rename

Rename

FILT\_HTN\_UNIQ\_ICD9\_COUNT

Joiner

Combine

Missing Value

Convert missing values to zero

# Machine learning and summing algorithms successfully replicate at Marshfield Clinic

	Model	Vanderbilt (n=631)			Replication Marshfield (n=100)		
		AUC	Sens.	PPV	AUC	Sens.	PPV
Random Forests	ICD9, meds, all BP (random forest)	0.955	0.844	0.954	0.922	0.966	0.919
	ICD9, meds, all vitals (random forest)	0.961	0.858	0.954	0.910	0.966	0.905
	<b>ICD9, meds, all vitals, RegEx (random forest)*</b>	0.967	0.866	0.954	<b>0.934</b>	0.966	0.934
	<b>ICD9, meds, all vitals, concept (random forest)</b>	<b>0.976</b>	0.902	0.952	0.873	0.966	0.864
	ICD9, meds, all vitals, RegEx, concepts (random forest)*	0.968	0.877	0.954	0.898	0.966	0.891

**\*Marshfield Clinic inputs to random forest models did not include regular expression (RegEx) information**

# Machine learning and summing algorithms successfully replicate at Marshfield Clinic

		Vanderbilt (n=631)			Replication Marshfield (n=100)		
Model		AUC	Sens.	PPV	AUC	Sens.	PPV
Random Forests	ICD9, meds, all BP (random forest)	0.955	0.844	0.954	0.922	0.966	0.919
	ICD9, meds, all vitals (random forest)	0.961	0.858	0.954	0.910	0.966	0.905
	<b>ICD9, meds, all vitals, RegEx (random forest)*</b>	0.967	0.866	0.954	<b>0.934</b>	0.966	0.934
	<b>ICD9, meds, all vitals, concept (random forest)</b>	<b>0.976</b>	0.902	0.952	0.873	0.966	0.864
	ICD9, meds, all vitals, RegEx, concepts (random forest)*	0.968	0.877	0.954	0.898	0.966	0.891
Category Counts	Count of HTN ICD9, med, and BP 2 of 3	0.833	0.952	0.822	0.646	1.000	0.670
	<b>Count of HTN ICD9, med, and BP 3 of 3</b>	0.877	0.798	0.967	<b>0.914</b>	0.949	0.918
	<b>Count of HTN ICD9, med, BP, and concept 3 of 4</b>	<b>0.910</b>	0.925	0.924	0.711	0.983	0.716

**\*Marshfield Clinic inputs to random forest models did not include regular expression (RegEx) information**

# Machine learning and summing algorithms successfully replicate at Marshfield Clinic

		Vanderbilt (n=631)			Replication Marshfield (n=100)		
Model		AUC	Sens.	PPV	AUC	Sens.	PPV
Random Forests	ICD9, meds, all BP (random forest)	0.955	0.844	0.954	0.922	0.966	0.919
	ICD9, meds, all vitals (random forest)	0.961	0.858	0.954	0.910	0.966	0.905
	<b>ICD9, meds, all vitals, RegEx (random forest)*</b>	0.967	0.866	0.954	<b>0.934</b>	0.966	0.934
	<b>ICD9, meds, all vitals, concept (random forest)</b>	<b>0.976</b>	0.902	0.952	0.873	0.966	0.864
	ICD9, meds, all vitals, RegEx, concepts (random forest)*	0.968	0.877	0.954	0.898	0.966	0.891
Category Counts	Count of HTN ICD9, med, and BP 2 of 3	0.833	0.952	0.822	0.646	1.000	0.670
	<b>Count of HTN ICD9, med, and BP 3 of 3</b>	0.877	0.798	0.967	<b>0.914</b>	0.949	0.918
	<b>Count of HTN ICD9, med, BP, and concept 3 of 4</b>	<b>0.910</b>	0.925	0.924	0.711	0.983	0.716
Category Sums	Sum of normalized HTN ICD9, meds, and BP	0.915	1.000	0.673	<b>0.949</b>	1.000	0.702
	<b>Sum of normalized HTN ICD9, meds, BP, and concept</b>	<b>0.929</b>	1.000	0.663	<b>0.949</b>	1.000	0.702

\*Marshfield Clinic inputs to random forest models did not include regular expression (RegEx) information



# Conclusions

- Combining categories improves performance
  - 0.976 AUC for best random forest
  - 0.908 AUC billing codes
  - 0.854 AUC hypertensive blood pressure readings
- Concepts are the most valuable individual category on Vanderbilt data
- All random forest models and summing algorithms performed comparably on Marshfield Clinic data