

Are Attribute Inference Attacks Just Imputation?

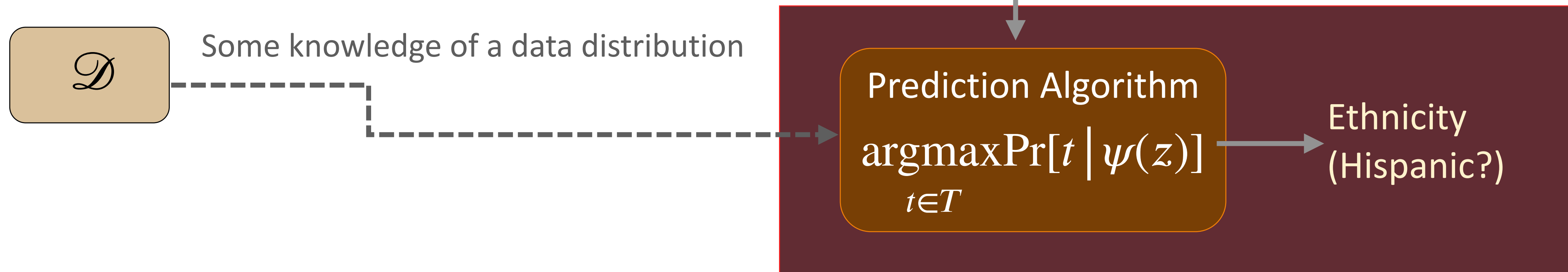
Bargav Jayaraman
bj4nq@virginia.edu
University of Virginia

David Evans
evans@virginia.edu
University of Virginia

Imputation

Infer missing attributes from available data

Hospital	Gender	Source	Stay Length	Patient Age	Ethnicity	Charges	...	Procedure
102	1	6	10	11	?	\$34920.33	...	34
102	0	2	3	8	?	\$4062.46	...	95
350	0	6	23	18	?	\$105239.23	...	62

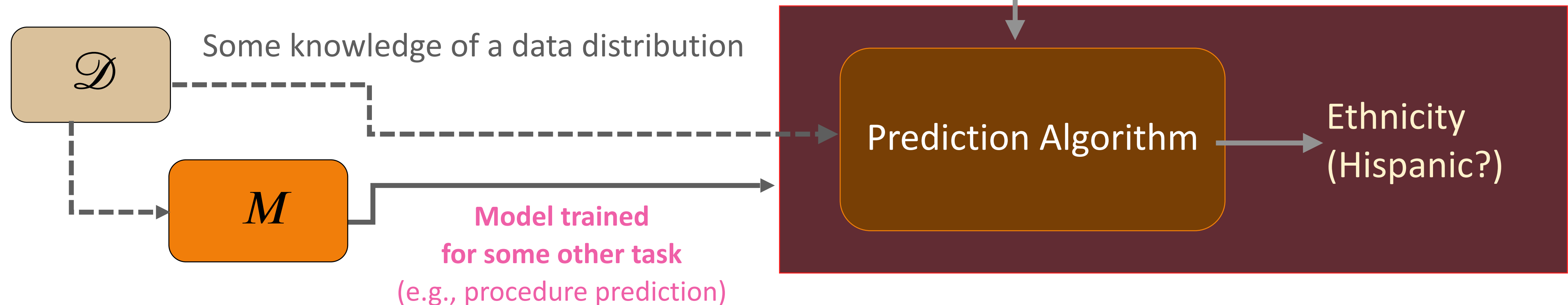


Traditionally, “imputation” is not considered an attack or privacy risk.

Attribute Inference

*Infer missing attributes from available data and **model***

Hospital	Gender	Source	Stay Length	Patient Age	Ethnicity	Charges	...	Procedure
102	1	6	10	11	?	\$34920.33	...	34
102	0	2	3	8	?	\$4062.46	...	95
350	0	6	23	18	?	\$105239.23	...	62



“Attribute inference” is considered an attack and a privacy risk.

Exploring Different Threat Settings

Adversary Knows the Training Distribution

Adversary's Data Set Size

Large

Yes

Prior AI Attacks

Fredrikson et al. [USENIX Sec 14]

Yeom et al. [CSF 18]

Mehnaz et al. [USENIX Sec 22]

Imputation
itself does well!

No

Small

Average Prediction Accuracy

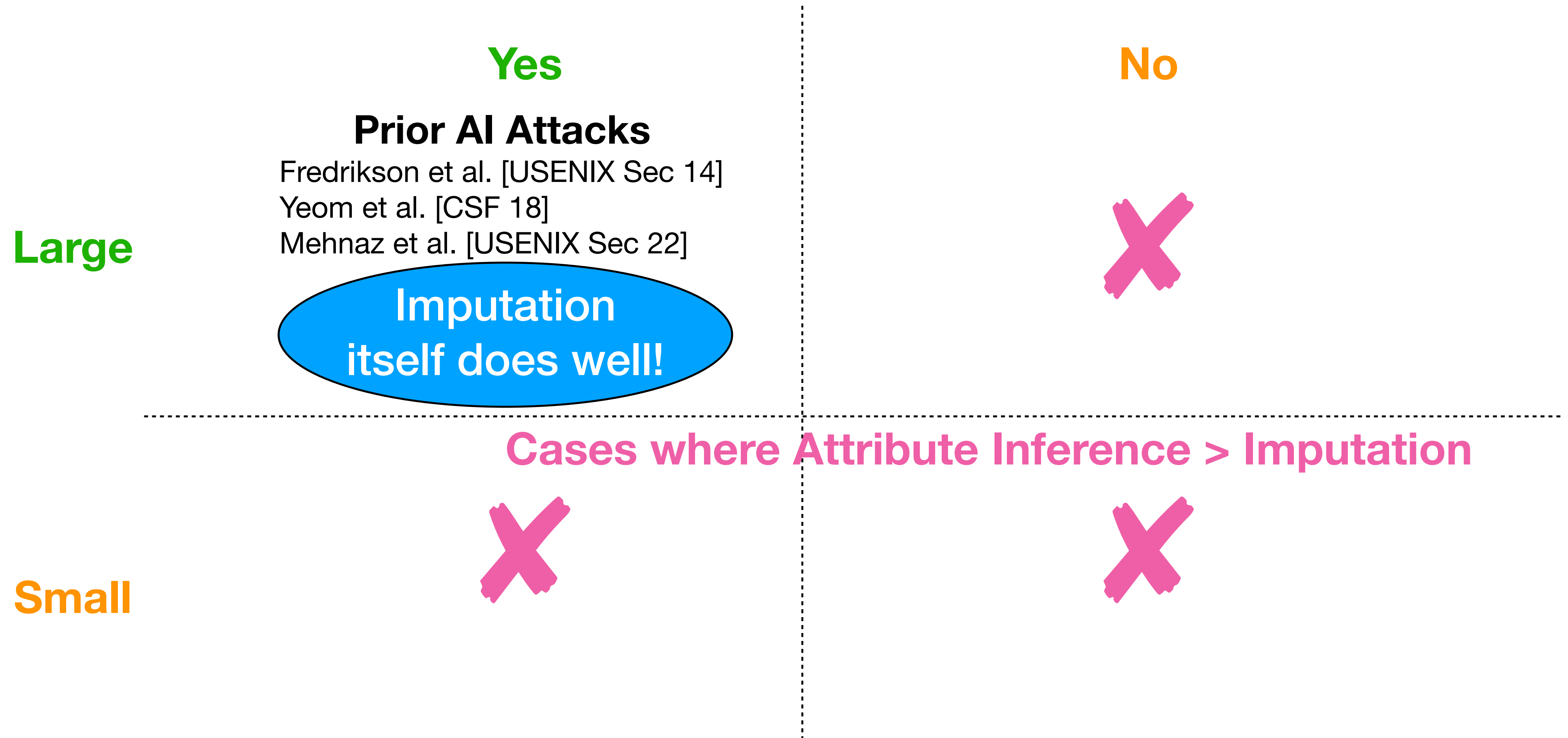
Attribute Prediction Method	Texas-100X		Census19	
	Gender	Ethnicity	Gender	Race
Naïve Most Common	0.62	0.72	0.52	0.78
Imputation	0.66	0.72	0.59	0.82
Yeom et al. [CSF 2018]	0.62	0.64	0.63	0.06
Mehnaz et al. [USENIX Sec 2022]	0.59	0.60	0.63	0.06
WCAI (our improvements to Yeom)	0.68	0.74	0.64	0.83

Black-box attribute inference attacks do not meaningfully outperform imputation:
no evidence that access to model helps

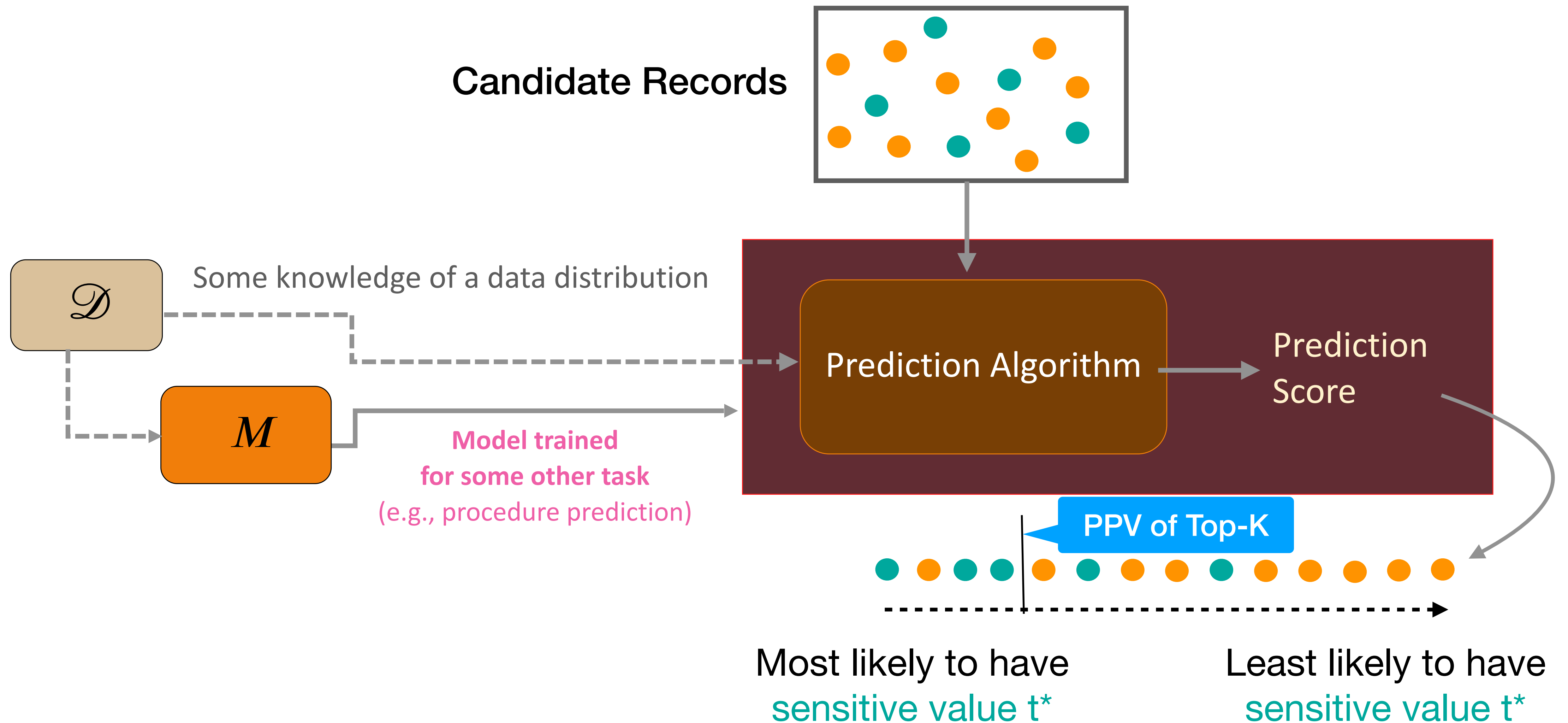
Exploring Different Threat Settings

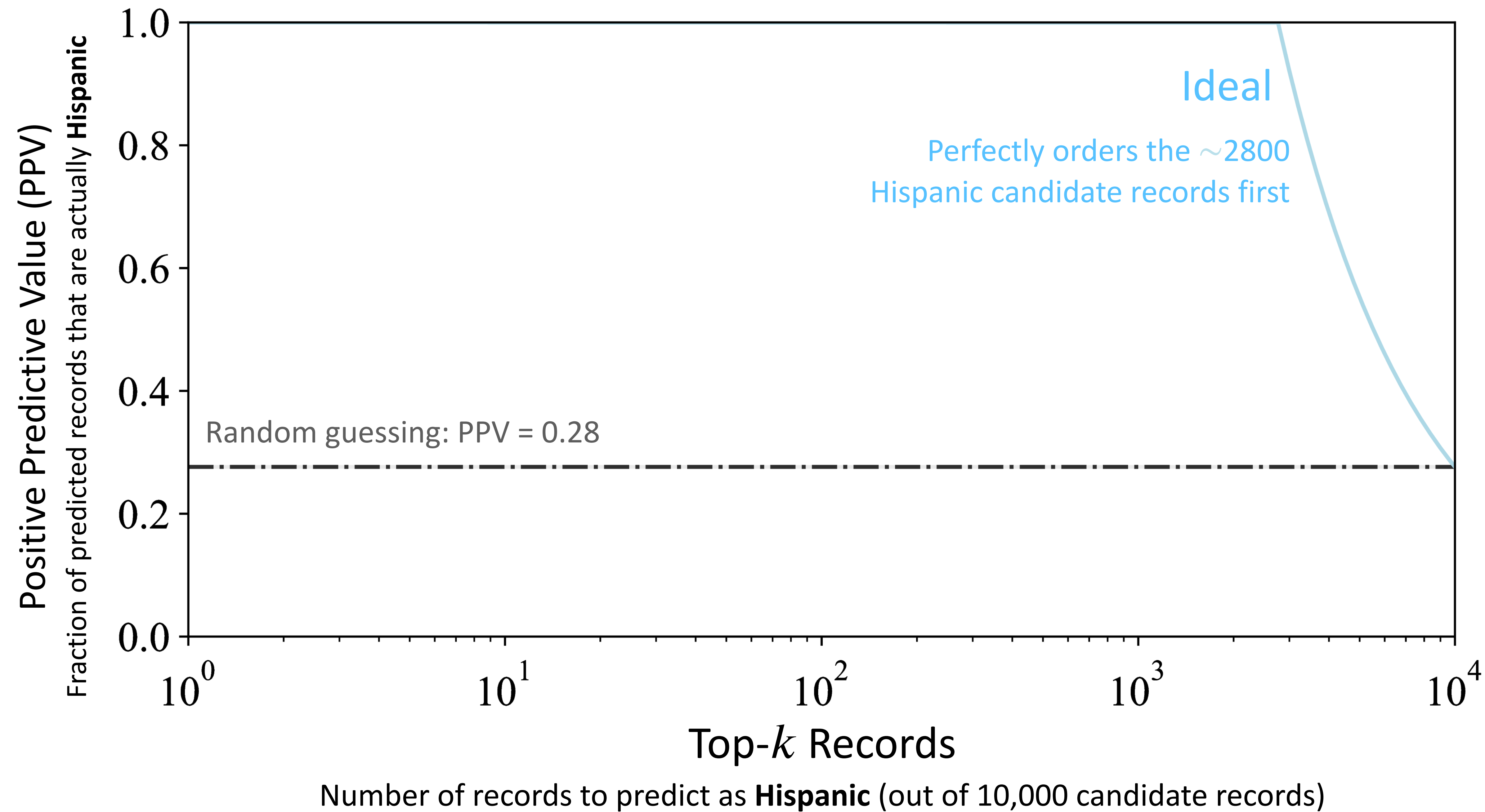
Adversary Knows the Training Distribution

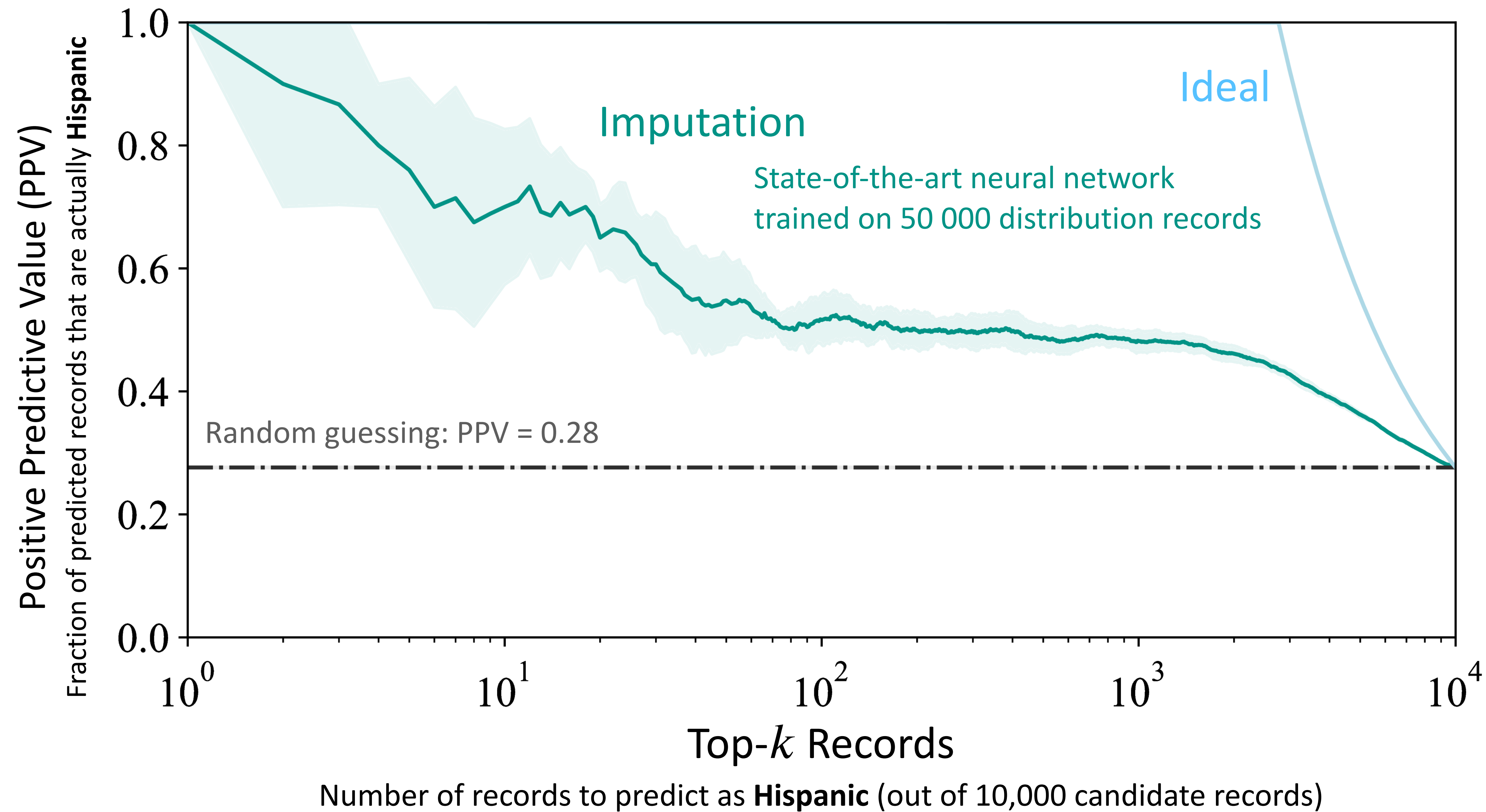
Adversary's Data Set Size

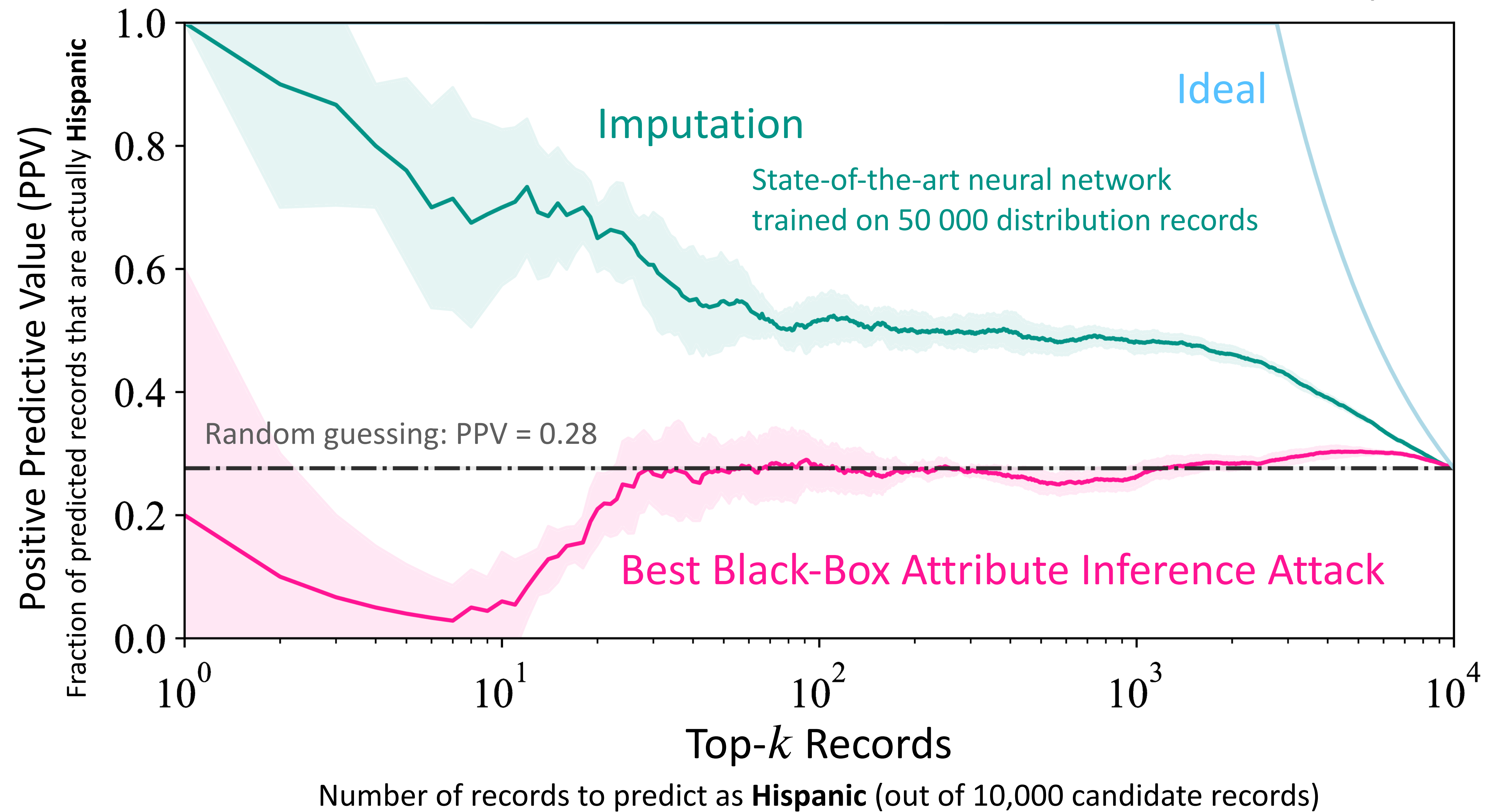


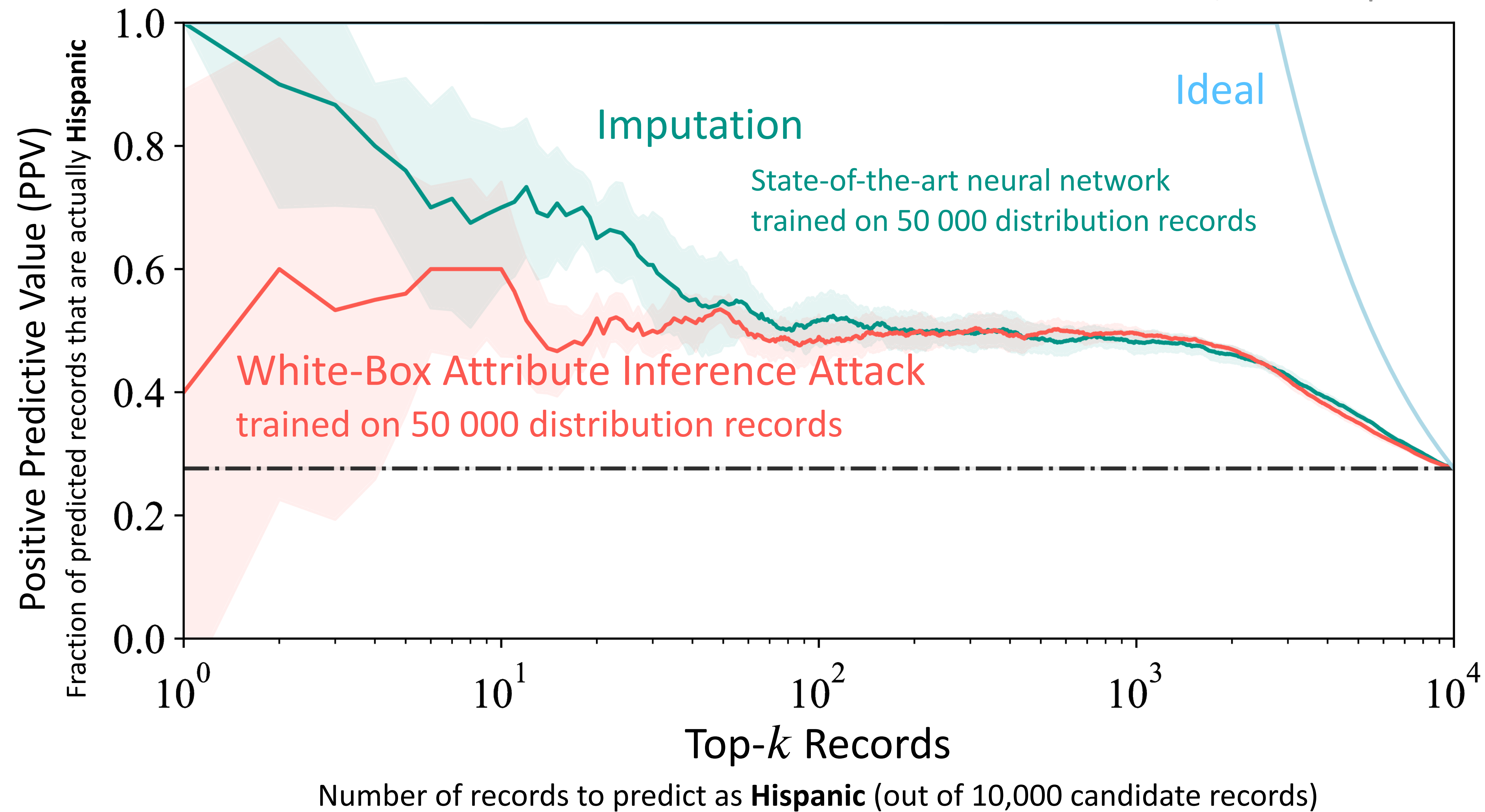
Sensitive Value Inference Attack



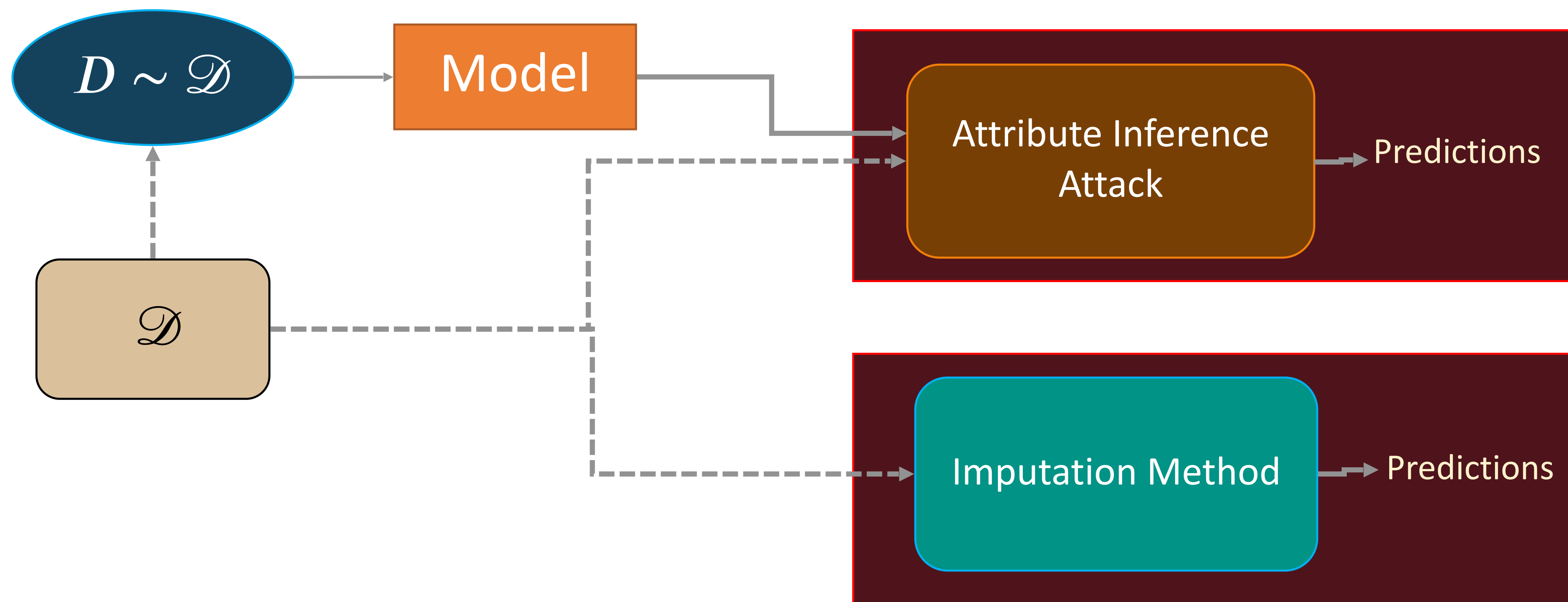






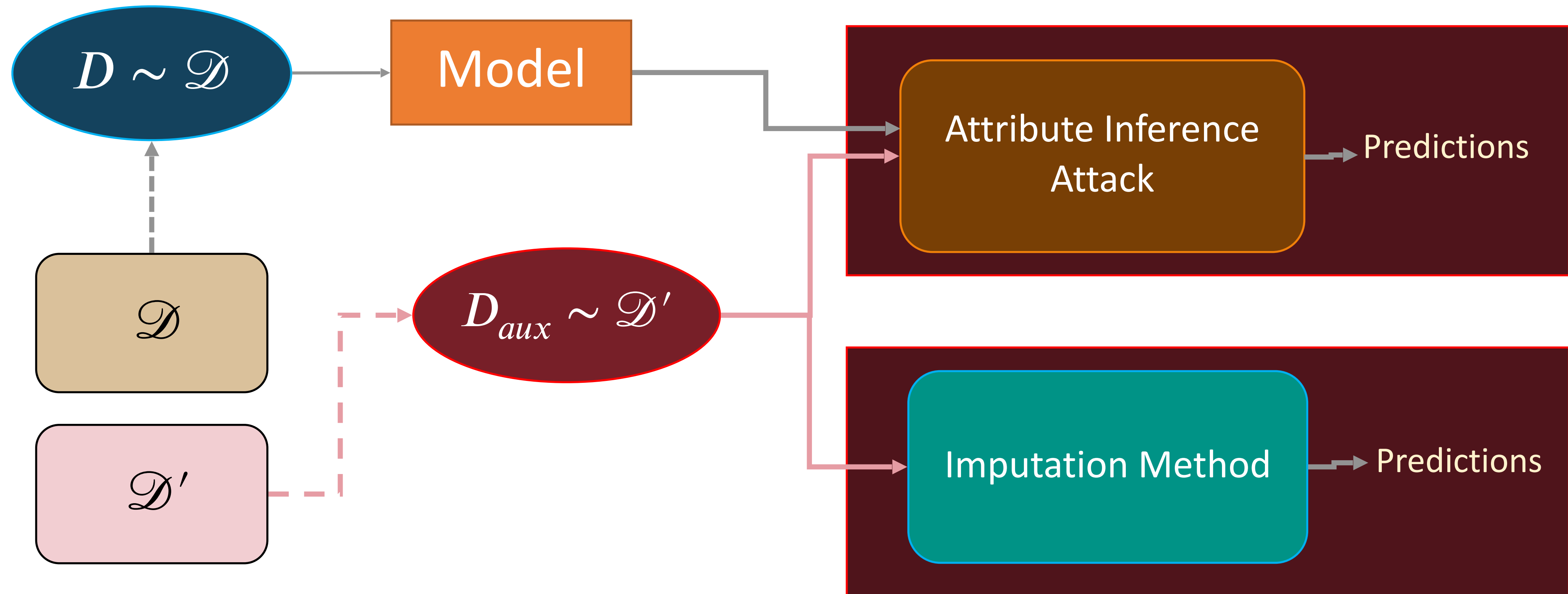


Are there Attribute Inference attacks that matter?

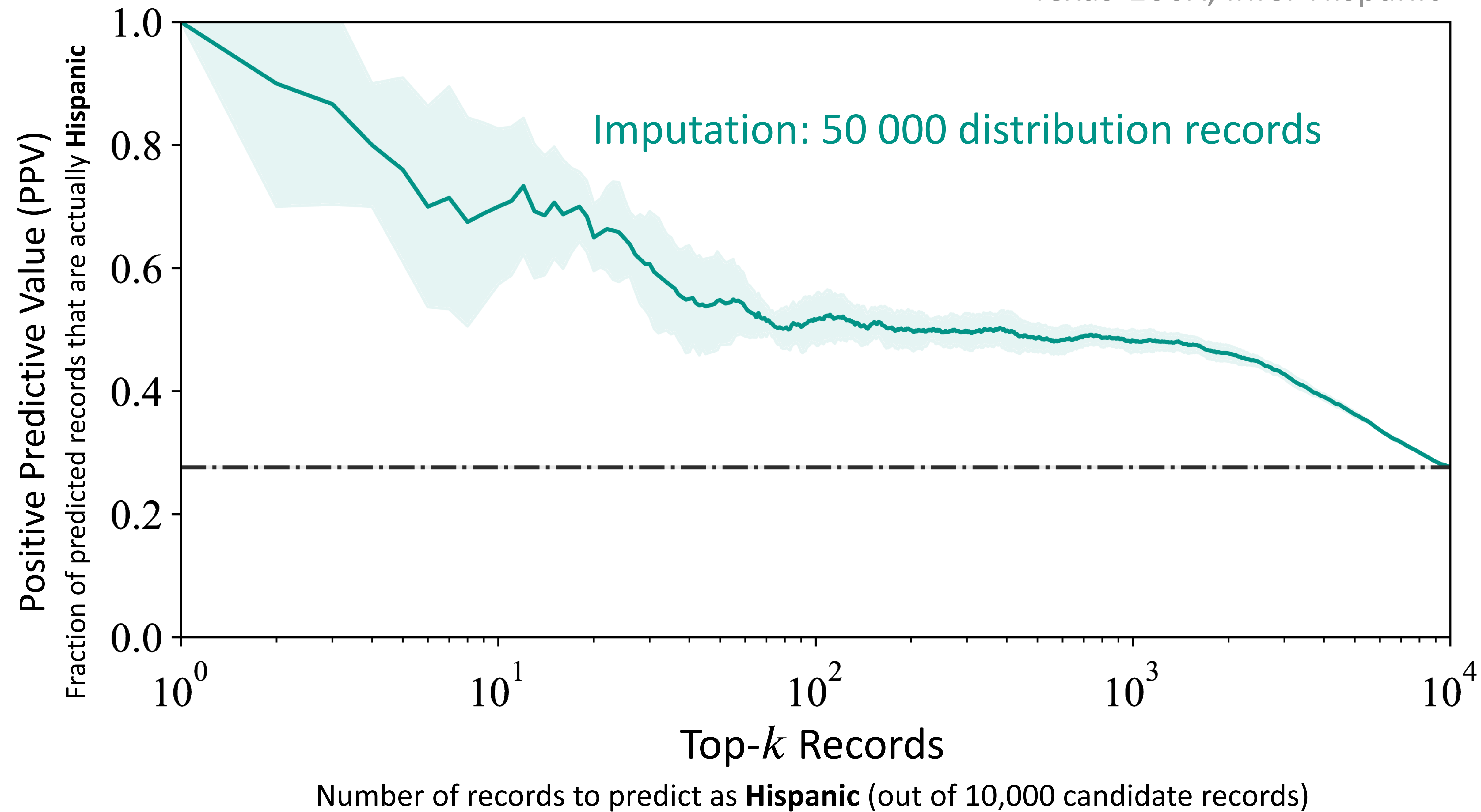


If adversary has good prior knowledge of training distribution \mathcal{D} ,
unlikely that model improves ability to infer attributes

Models Revealing Useful Information

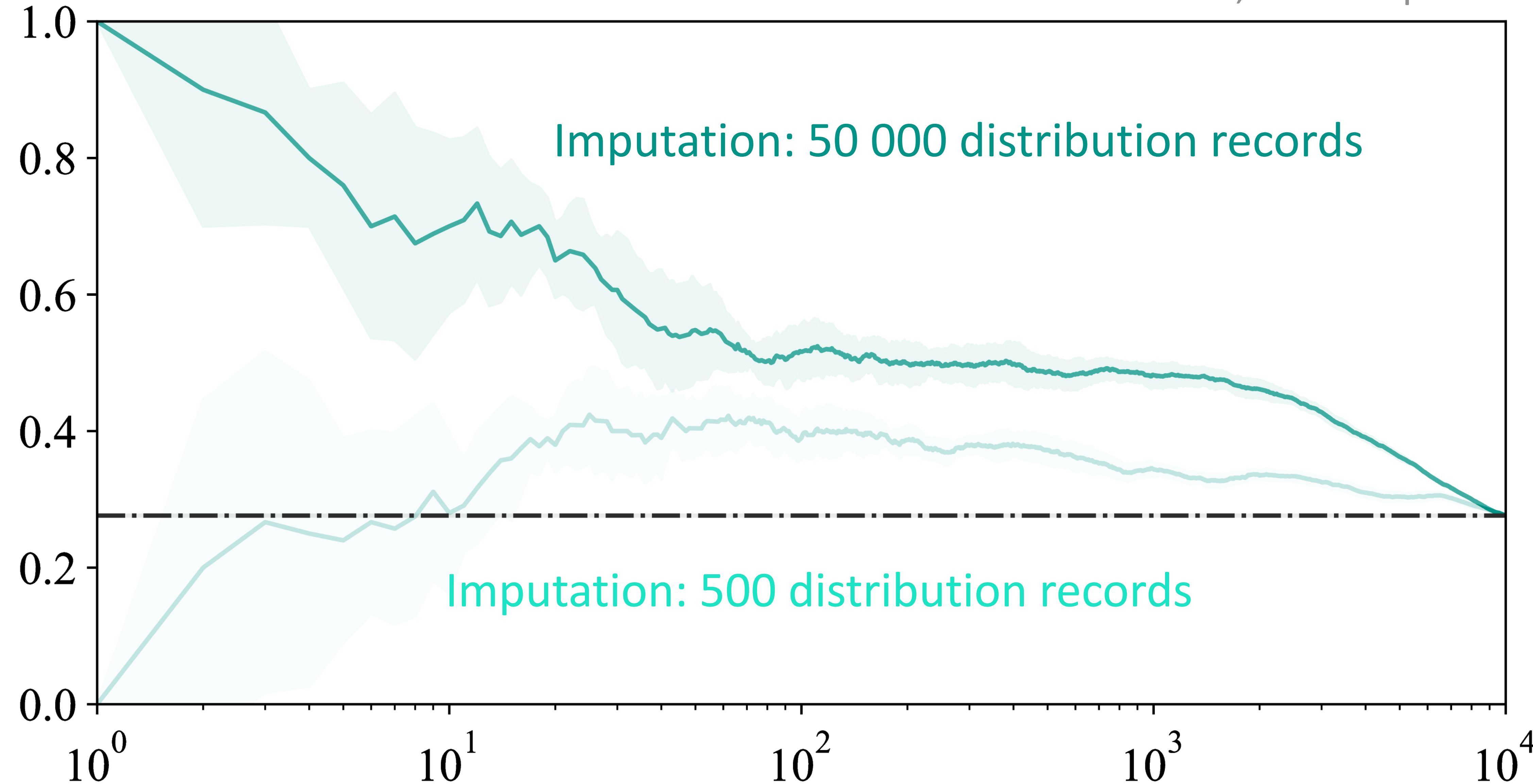


If adversary has *limited* prior knowledge of training distribution \mathcal{D} ,
model may improve ability to infer attributes



Positive Predictive Value (PPV)

Fraction of predicted records that are actually **Hispanic**



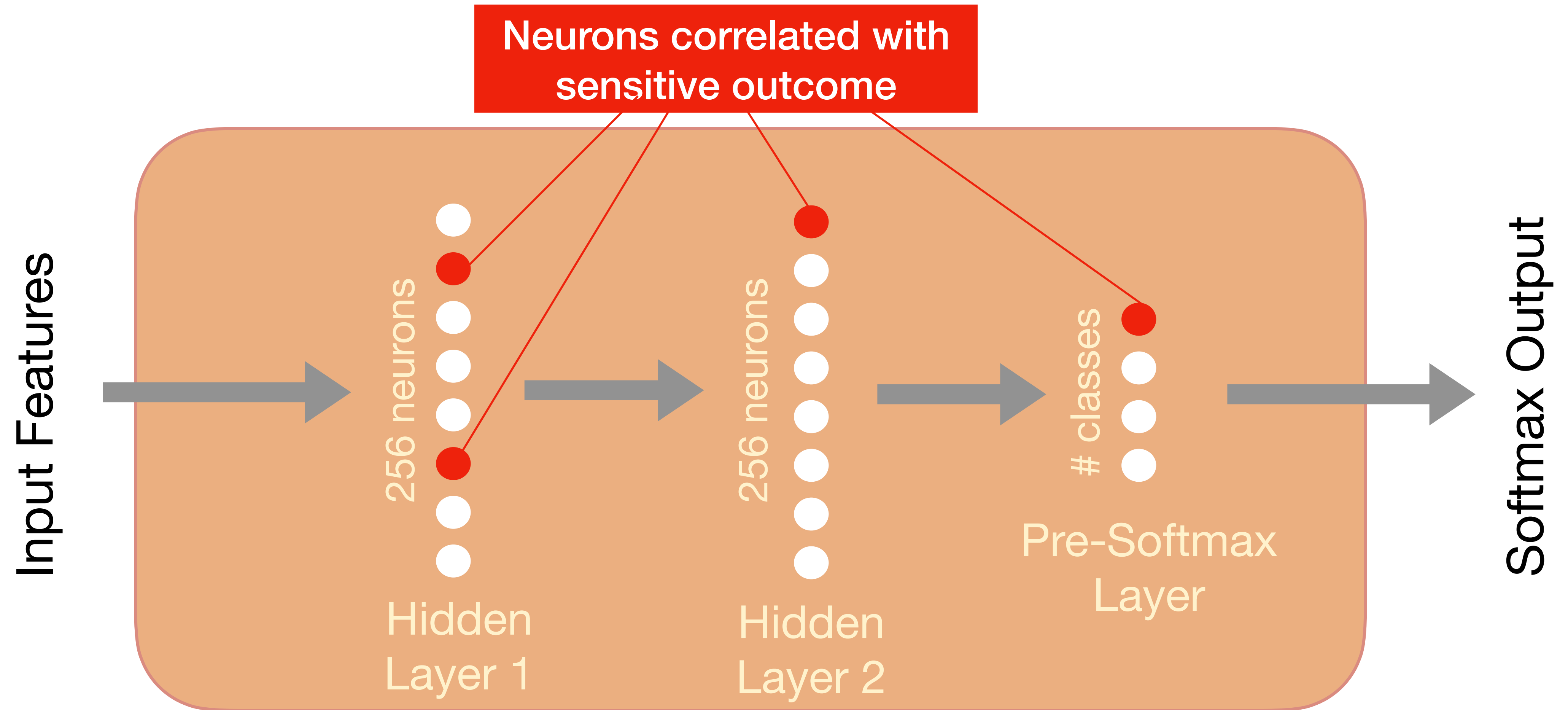
Imputation: 50 000 distribution records

Imputation: 500 distribution records

Top-*k* Records

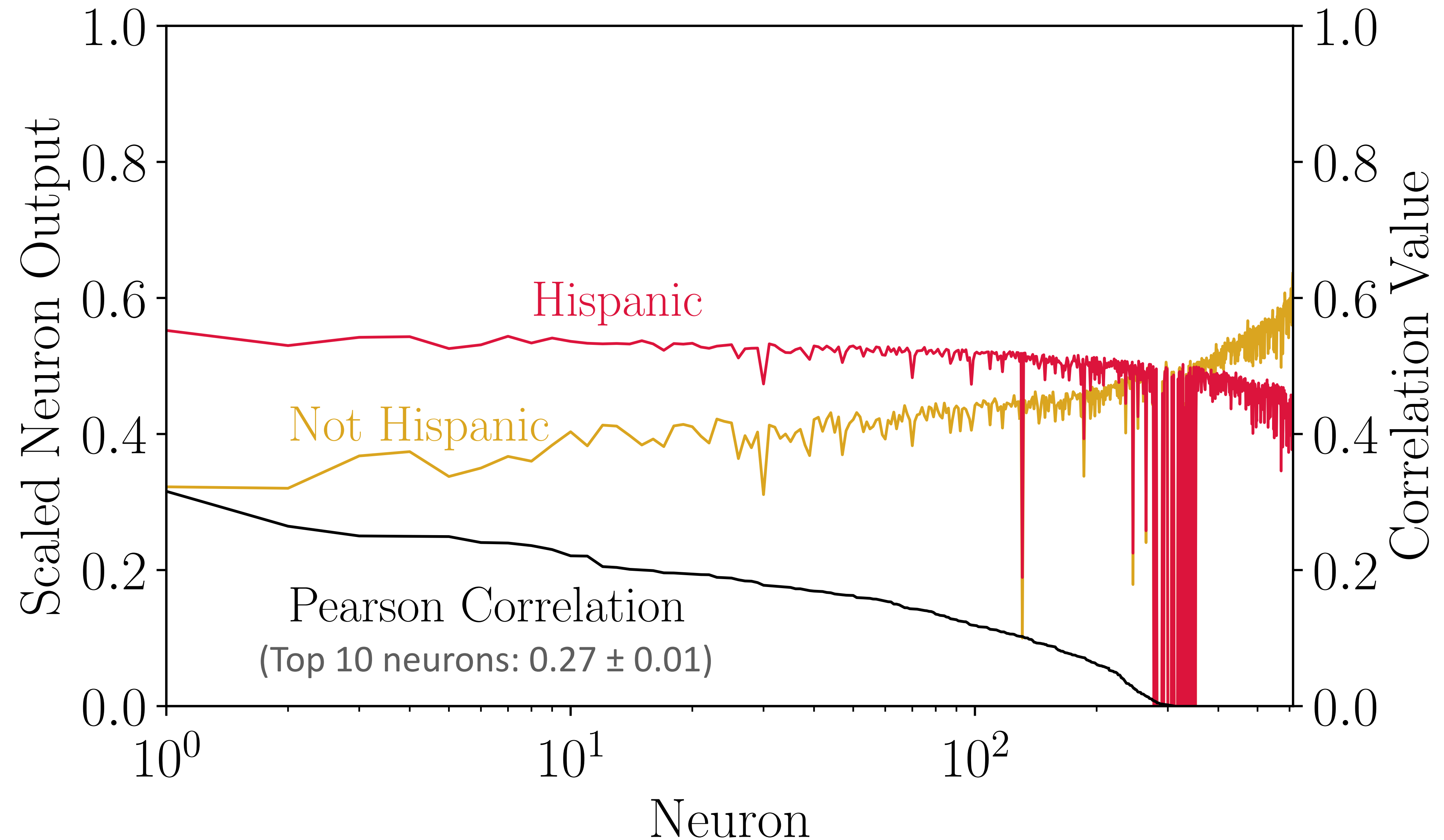
Number of records to predict as **Hispanic** (out of 10,000 candidate records)

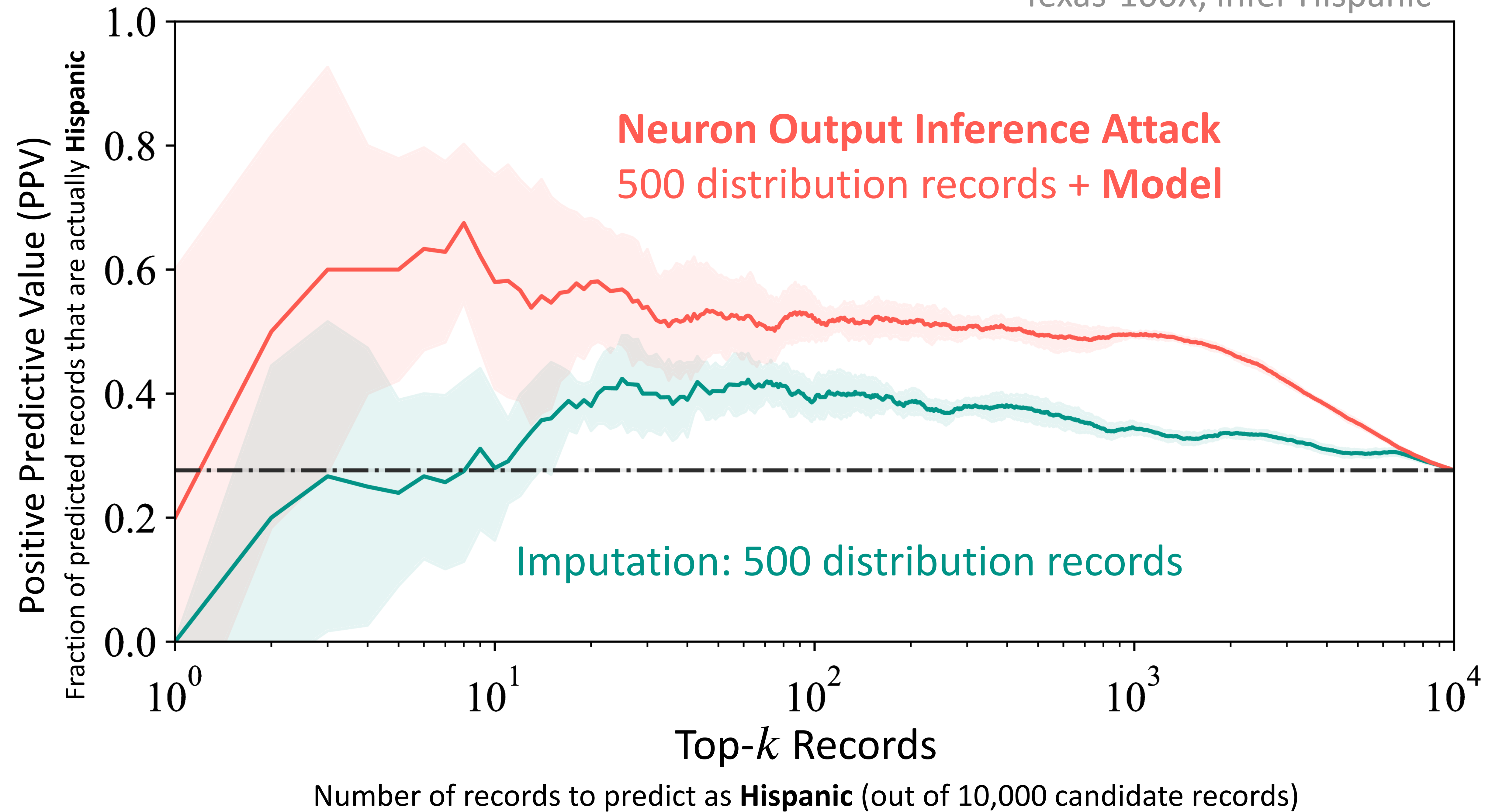
Neuron Output Inference Attack

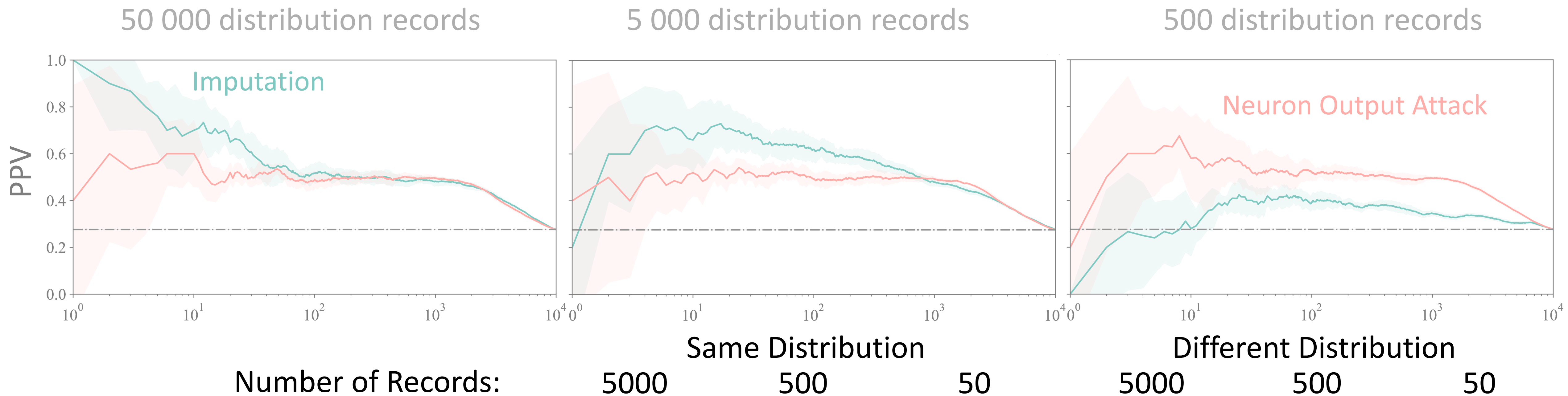


Neural Network Model Architecture

Some Neurons Activate Differently







		Same Distribution			Different Distribution		
		5000	500	50	5000	500	50
Texas-100X (Hispanic)	Imputation	0.62 ± 0.05	0.39 ± 0.03	0.24 ± 0.01	0.44 ± 0.02	0.41 ± 0.05	0.37 ± 0.05
	Neuron Output Attack	0.49 ± 0.02	0.52 ± 0.03	0.47 ± 0.05	0.49 ± 0.02	0.49 ± 0.04	0.52 ± 0.07
Census19 (Asian)	Imputation	0.91 ± 0.03	0.25 ± 0.02	0.55 ± 0.06	0.90 ± 0.03	0.46 ± 0.04	0.42 ± 0.04
	Neuron Output Attack	0.85 ± 0.03	0.82 ± 0.05	0.67 ± 0.06	0.86 ± 0.05	0.85 ± 0.04	0.65 ± 0.12

PPV at $k=100$ (1%) (highlighted when model provides significant benefit)

Differential Privacy does not Mitigate Attribute Inference Risk

Texas-100X, infer Hispanic

	Without DP	With DP	Training	Non-Training
Imputation	0.62 ± 0.05	0.62 ± 0.05	0.62 ± 0.05	0.63 ± 0.02
Neuron Attack	0.49 ± 0.02	0.49 ± 0.03	0.49 ± 0.03	0.48 ± 0.02

Results for models trained with $(\epsilon = 1, \delta = 10^{-5})$ –DP

No significant differences between non-DP models, and no differences between predictions for candidates from training and non-training data

Conclusion

Attribute Inference Attacks *are doing* Imputation

Privacy Risk when the Distribution is not Public

Presenter:

Bargav Jayaraman
bj4nq@virginia.edu
University of Virginia

Code Repository:

<https://github.com/bargavj/EvaluatingDPML>