

# Predicting Chromatography-tandem Mass Spectrometry Amenability to Improve Non-targeted Analysis

Charles N. Lowe<sup>1</sup>, Kristin Isaacs<sup>1</sup>, Andrew McEachran<sup>2</sup>, Chris Grulke<sup>1</sup>, Jon Sobus<sup>1</sup>, Elin Ulrich<sup>1</sup>, Alex Chao<sup>1</sup>, and Antony J. Williams<sup>1</sup>

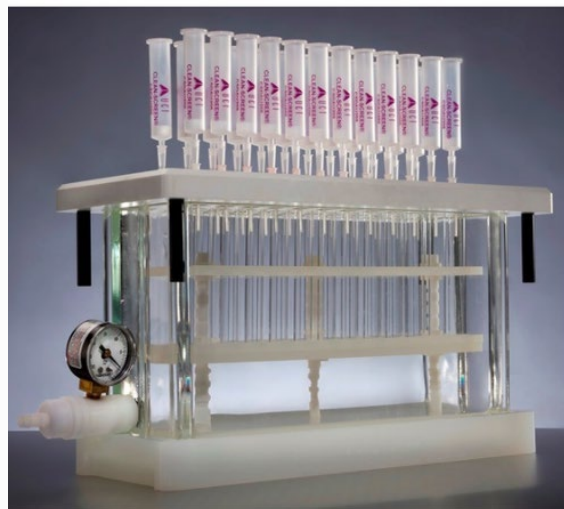
1. Center for Computational Toxicology and Exposure, U.S. Environmental Protection Agency, Research Triangle Park, NC
2. Agilent Technologies, Inc., Santa Clara, CA

**Disclaimer:** The views expressed in this presentation are those of the authors and do not necessarily reflect the views or policies of the U.S. Environmental Protection Agency.

# Complex samples, NTA, and the modeling problem



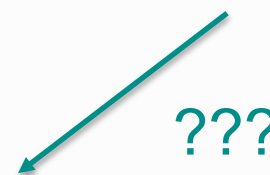
**Media Sample**



**Extraction, Cleanup &  
Sample Preparation**

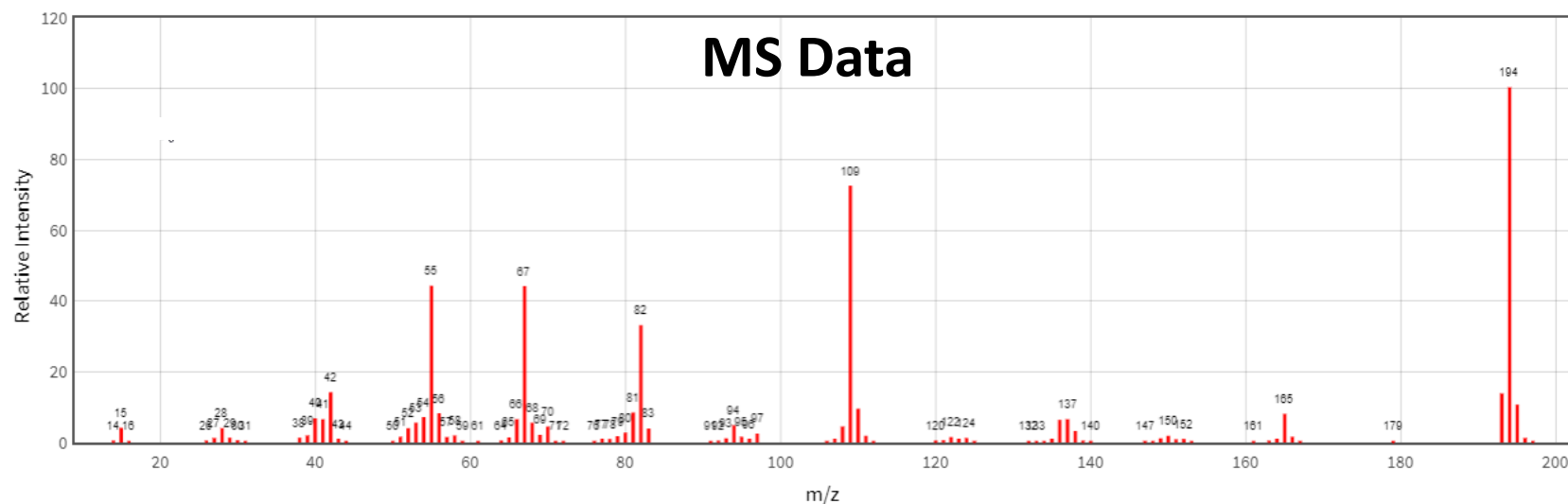


**MS Analysis**



Mass Spectrum

**MS Data**



# Curating a dataset for modeling


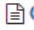
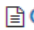
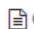

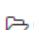
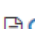
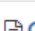
**MoNA - MassBank of North America** | Spectra | Downloads | Upload | Help

Search...

## Downloads

A set of commonly referenced predefined queries. Clicking the name of the query will display the associated spectra in the query browser. Each query is also available to download in either the MoNA internal JSON format or as NIST MS Search compatible MSP files.

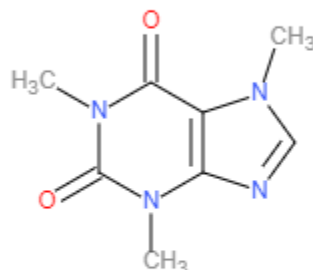
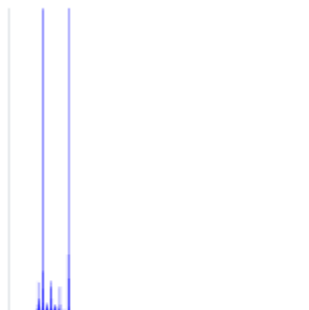
☐ Display Hidden Downloads

 <a href="#">Q All Spectra (659,728 spectra)</a>	<a href="#">Download</a>
 <a href="#">Q In-Silico Spectra (490,087 spectra)</a>	<a href="#">Download</a>
 <a href="#">Q Experimental Spectra (169,641 spectra)</a>	<a href="#">Download</a>
 <a href="#">Q GC-MS Spectra (18,883 spectra)</a>	<a href="#">Download</a>
 <a href="#">Q LC-MS Spectra (133,301 spectra)</a>	<a href="#">Download</a>
 <a href="#">Q LC-MS/MS Spectra (125,833 spectra)</a>	<a href="#">Download</a>
 <a href="#">Q LC-MS/MS Positive Mode (86,576 spectra)</a>	<a href="#">Download</a>
 <a href="#">Q LC-MS/MS Negative Mode (38,475 spectra)</a>	<a href="#">Download</a>

- 772 compounds in derivatized GCMS
- 7,199 compounds in non-derivatized GCMS
- **4,145 compounds in ESI+ LCMS**
- **2,981 compounds in ESI- LCMS**

## Caffeine

Score: ★★★★★

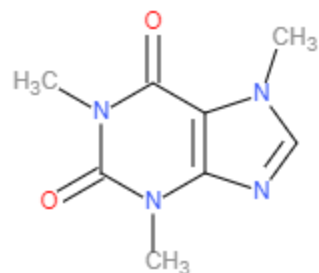


instrument	Pegasus III TOF-MS system...
instrument type	GC-EL-TOF
ms level	MS1
retention index	1880.2430
retention time	724.344 sec
ionization mode	positive
accession	OUF00133
date	2016.01.19 (Created 2010....
author	Tsujimoto Y, Tsugawa H, B...
license	CC BY-SA

Originally submitted to the [MassBank High Quality Mass Spectral Database](#)

## Caffeine

Score: ★★★★★



instrument type	QqQ
instrument	Micromass Quattromicro
collision energy	15eV
ionization	ESI
ionization mode	positive
ms level	MS2
precursor m/z	194.9000
precursor type	[M+H] <sup>+</sup>
accession	PM018511
publication	Alonso-Salces RM, Guillou...

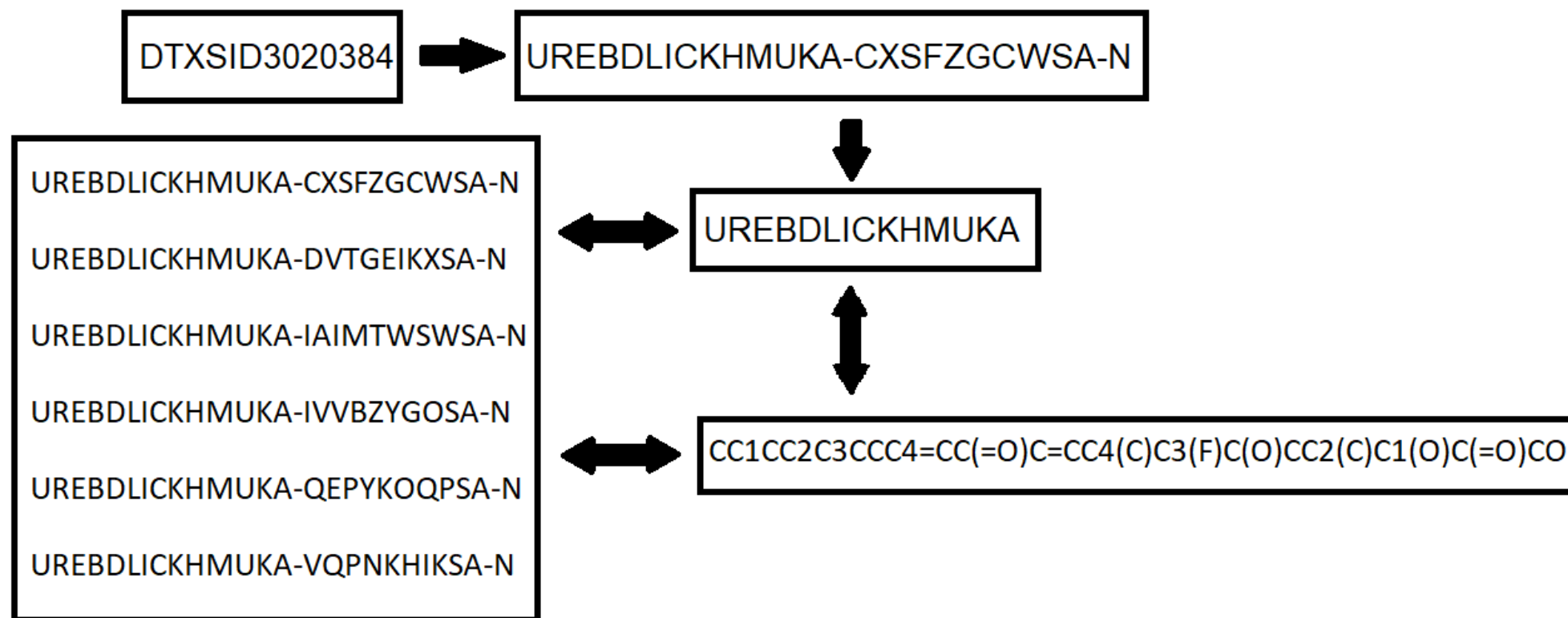
Originally submitted to the [RIKEN MS<sup>n</sup> Spectral Database for Phytochemicals](#)

# Curating a dataset for modeling

- Only amenable compounds identified in MoNA
  - No unamenable compound data
- ToxCast library LCMS curation
  - Spectra checked individually for quality
    - Provides unamenable compound data
- ESI+ LCMS
  - 403 amenable; 469 unamenable
- ESI- LCMS
  - 464 amenable; 415 unamenable
- Caveat: some of these unamenable compounds are amenable based on MoNA\*



# Curating a dataset for modeling



# Describing molecular structures

**Software News and Update**  
**PaDEL-Descriptor: An Open Source Software to  
Calculate Molecular Descriptors and Fingerprints**

**CHUN WEI YAP**

*Department of Pharmacy, Pharmaceutical Data Exploration Laboratory,  
National University of Singapore, Singapore*

*Received 17 May 2010; Revised 22 August 2010; Accepted 12 October 2010*

*DOI 10.1002/jcc.21707*

*Published online 17 December 2010 in Wiley Online Library (wileyonlinelibrary.com).*

- 1,444 1D & 2D Molecular descriptors from QSAR-ready SMILES. Examples include...
  - Electrotological states weighted by atomic properties
  - molecular linear free energy relationships weighted by atomic properties
  - Atom, bond, & ring counts
  - logP predictions, etc..



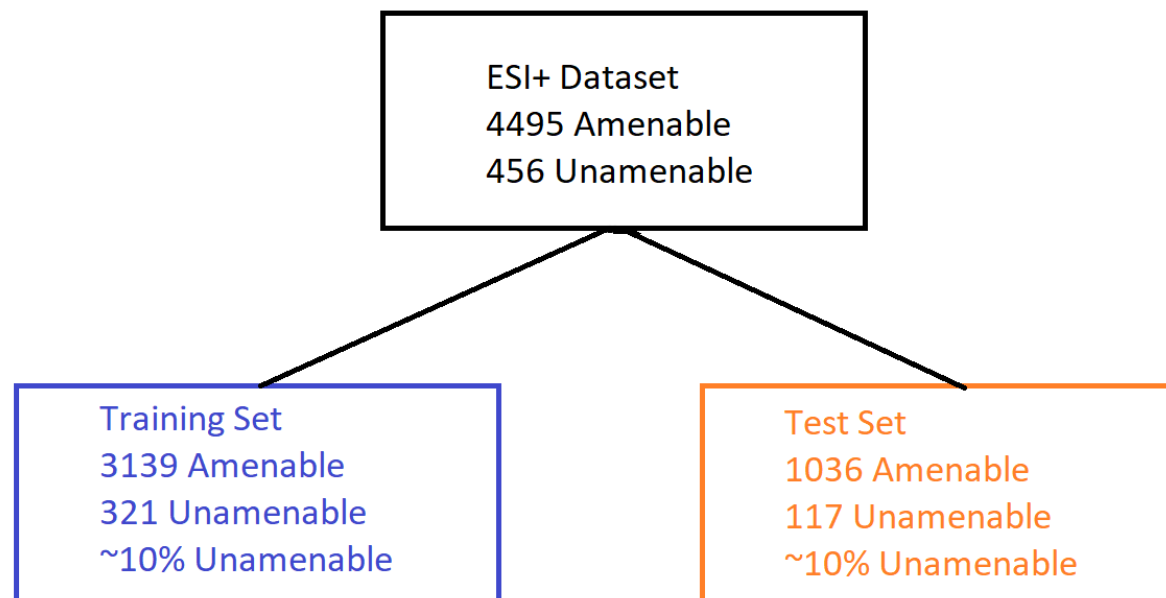
# Cleaning and reduction of descriptor space

- Dimension reduction will do two things:
  - improve interpretability of models
  - make model calculations faster
- Remove chemicals missing descriptors\*
- Remove any constant descriptors ( $\text{variance}(x) = 0$ )
- Remove near-constant descriptors ( $\text{sd}(x) < 0.25$ )
  - 0.25 gives a good balance between reduction and retention
- Calculate pair-wise correlations between remaining descriptors
  - Eliminate based on a cutoff = 0.96 correlation
    - descriptor showing largest pair correlation with other descriptors was excluded

**1,444 descriptors → 498 descriptors**

# Datasets suitable for modeling

- Models randomly divided into training and test sets
  - 75% of data for training, 25% for testing
  - Data stratified to maintain proportions in outcome variable
  - Different for each model
  - InChIKey skeleton as identifier



# Machine learning approach

- Random forest models for two endpoints
  - ESI+ LCMS, ESI- LCMS
  - Balance training set with either upsampling or downsampling
  - Optimize hyperparameters via grid search
    - Number of decision trees
    - Number of random descriptors selected at each node
  - 5-fold cross validation
  - Y-randomization
    - Randomly scramble endpoint, descriptors left intact

## Performance metrics

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \frac{TN}{TN + FP}$$

$$\text{balanced accuracy} = \frac{\text{sensitivity} + \text{specificity}}{2}$$

# Machine learning approach

## Random Forest Algorithm

Training set  $X = x_1 x_2 \dots x_n$  with responses

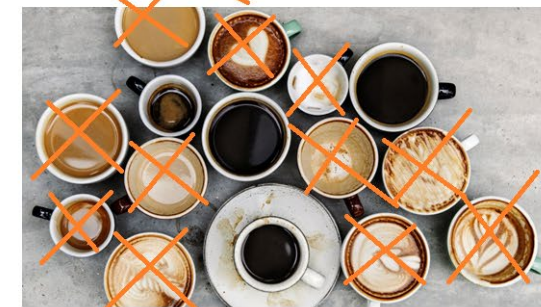
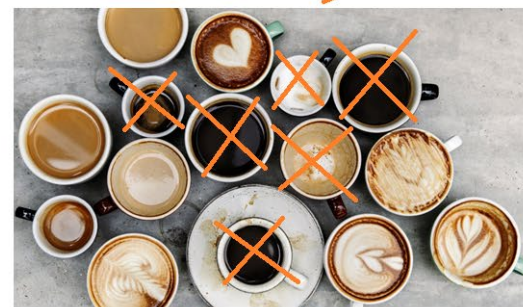
$$Y = y_1 y_2 \dots y_n$$

For  $b = 1, \dots, B$

1. Sample, with replacement,  $n$  training examples from  $X, Y$ ;  $X_b, Y_b$ .
2. Train a classification tree  $f_b$  on  $X_b, Y_b$ .
3. The majority of all  $f_b$  classifies unseen endpoints.



Yes      Creamer?      No

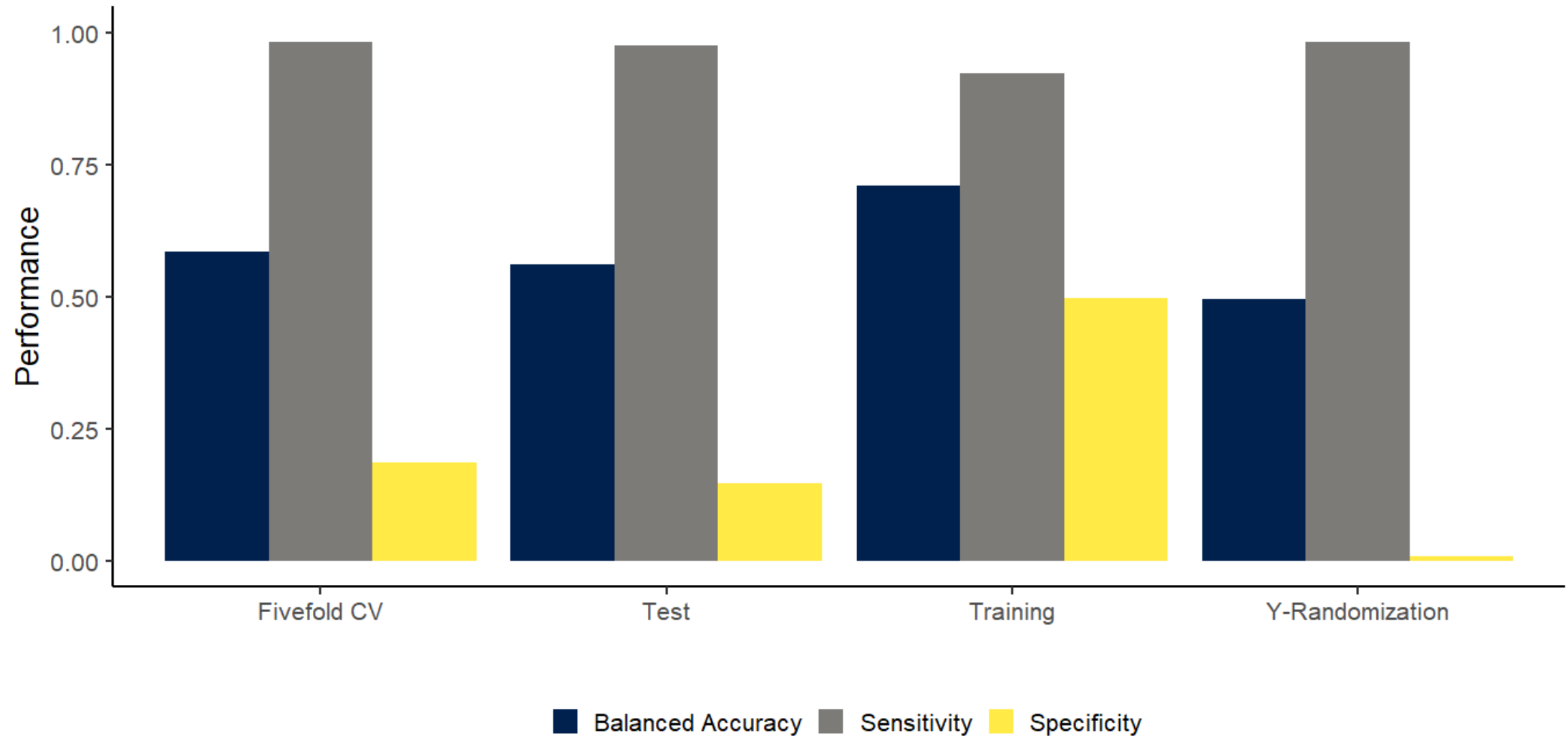


Yes      Artwork?      No



# Model performance

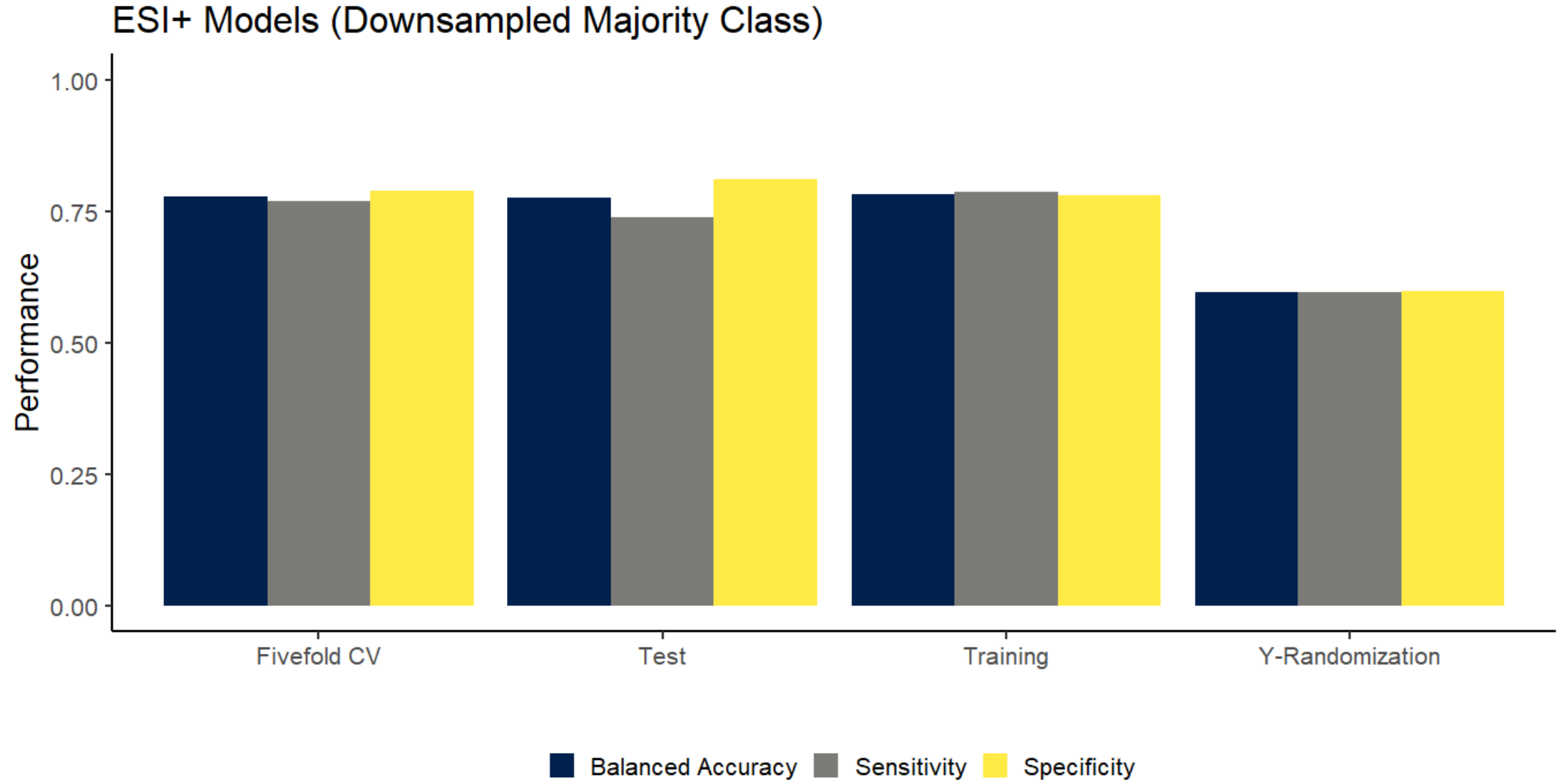
ESI+ Models (No Subsampling Applied)



# Model summary

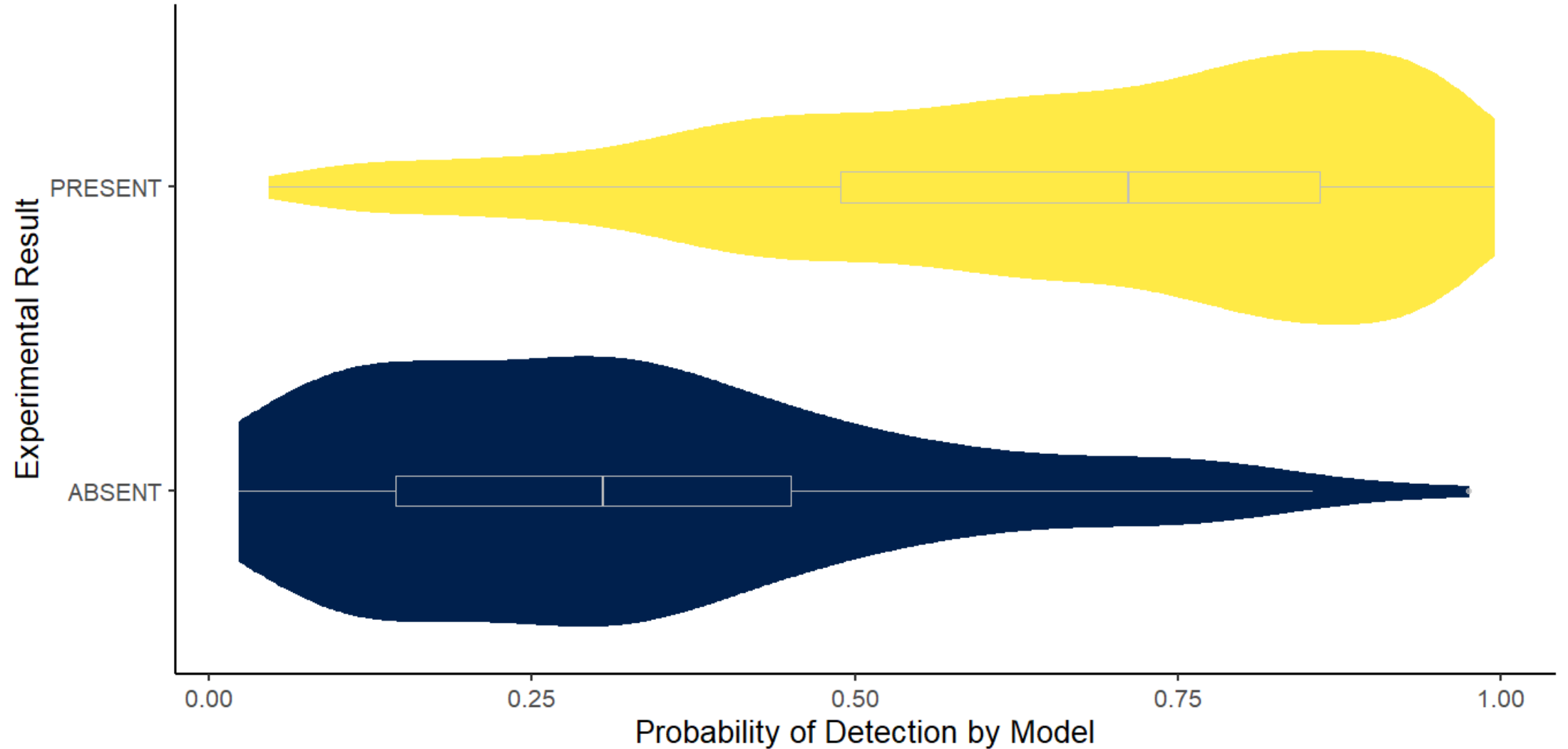
Imbalanced training sets lead to bad models!

# Model performance



# Model performance

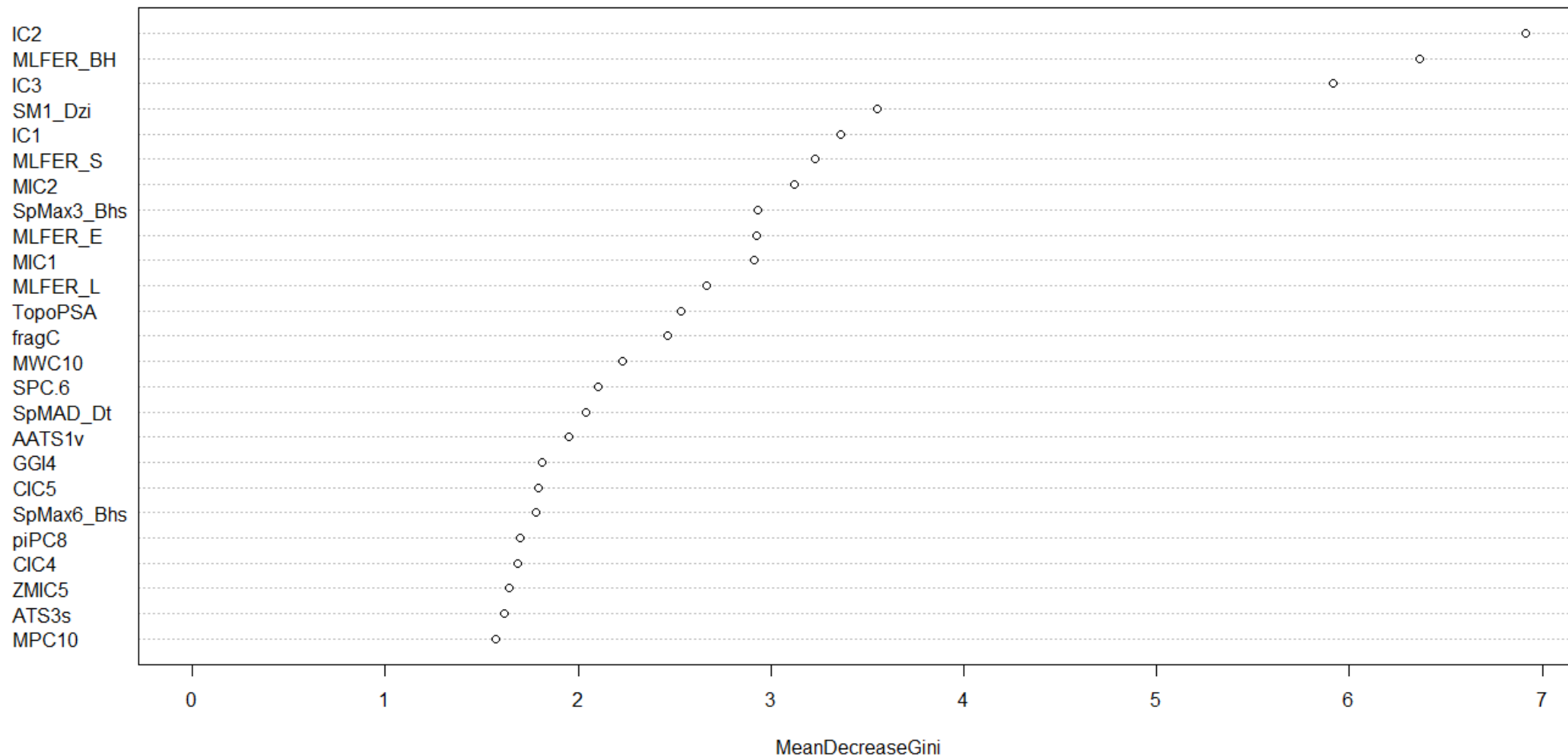
ESI+ Models (Downsampled Majority Class)





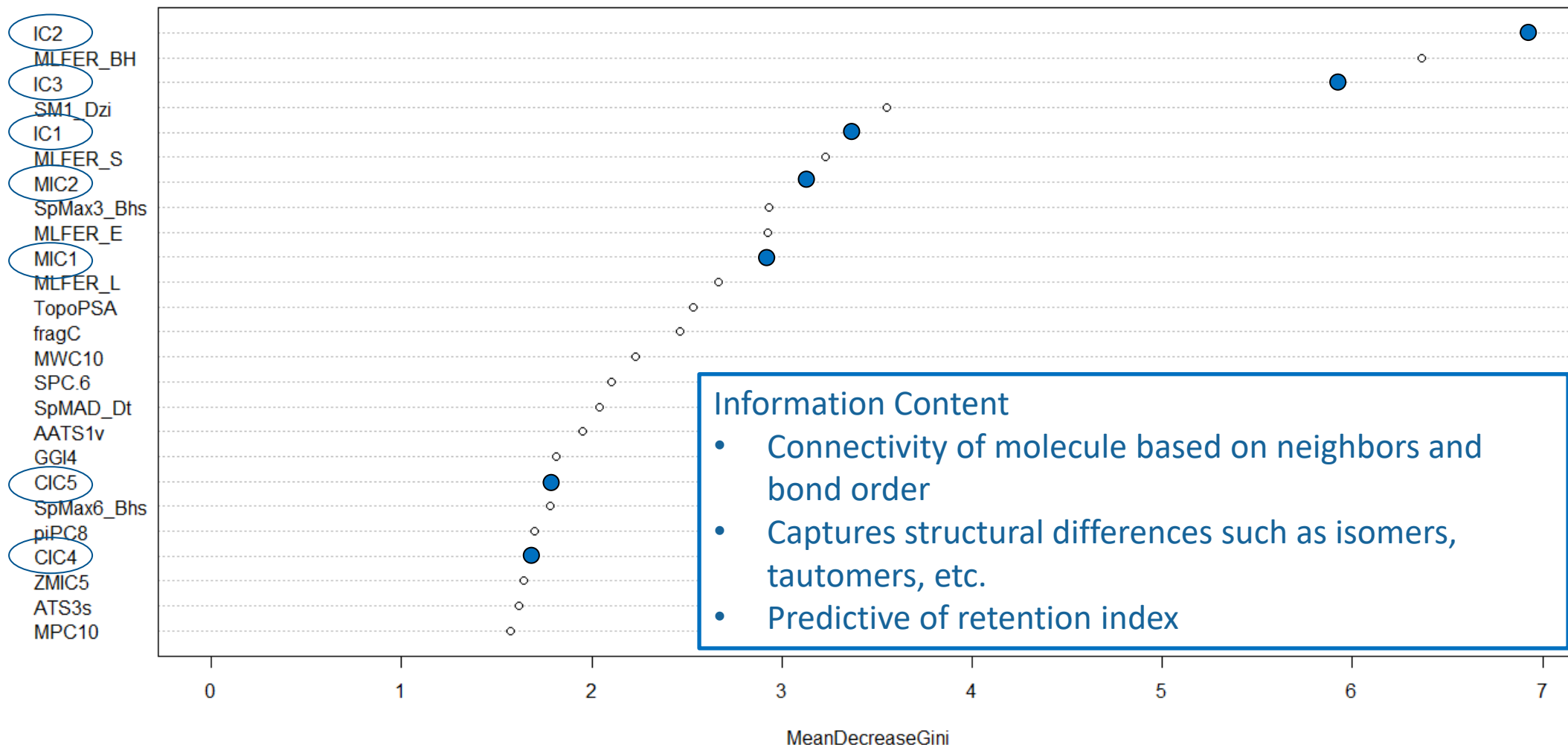
# Model performance

**ESI+ Models (Downsampled Majority Class)**



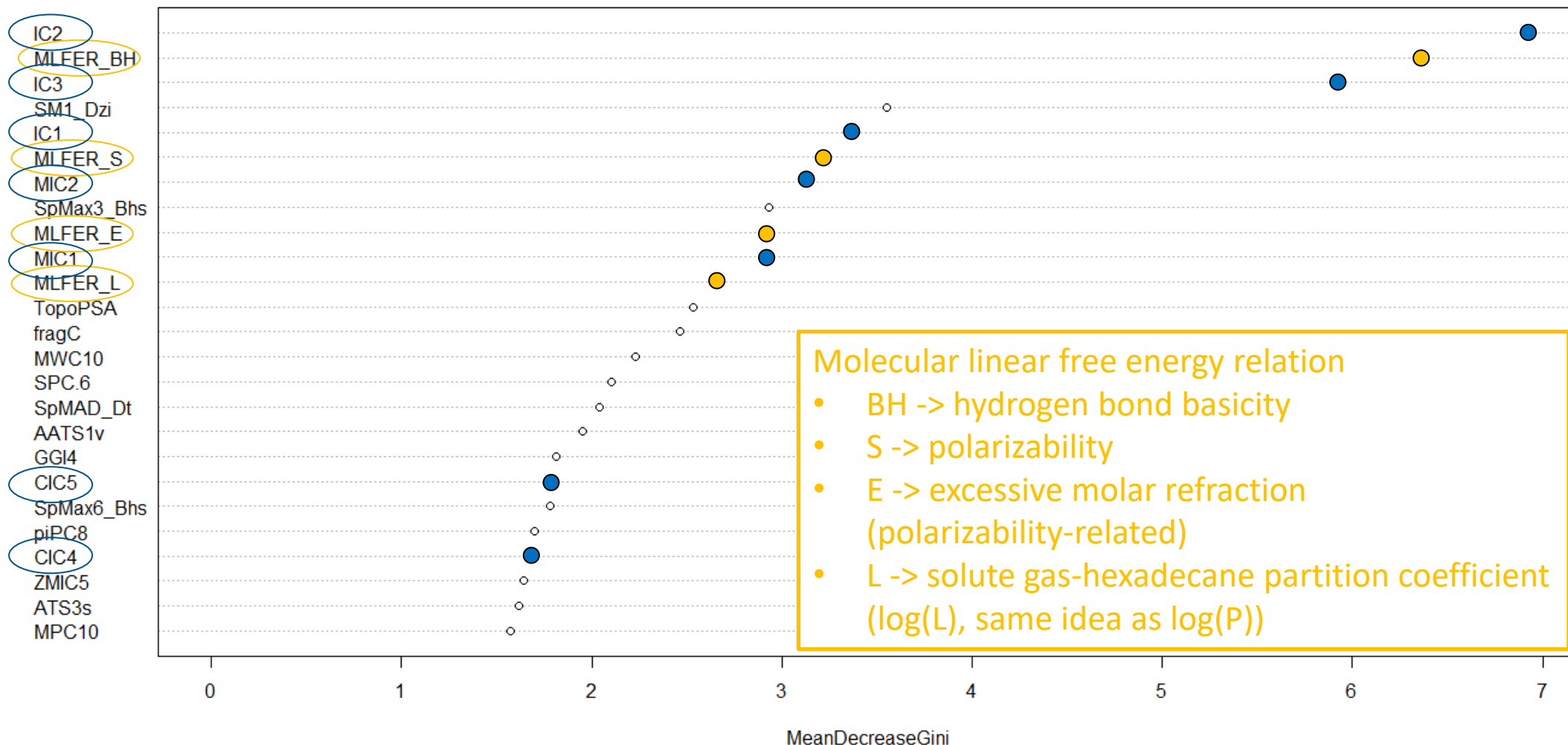
# Mechanistic interpretation

ESI+ Models (Downsampled Majority Class)



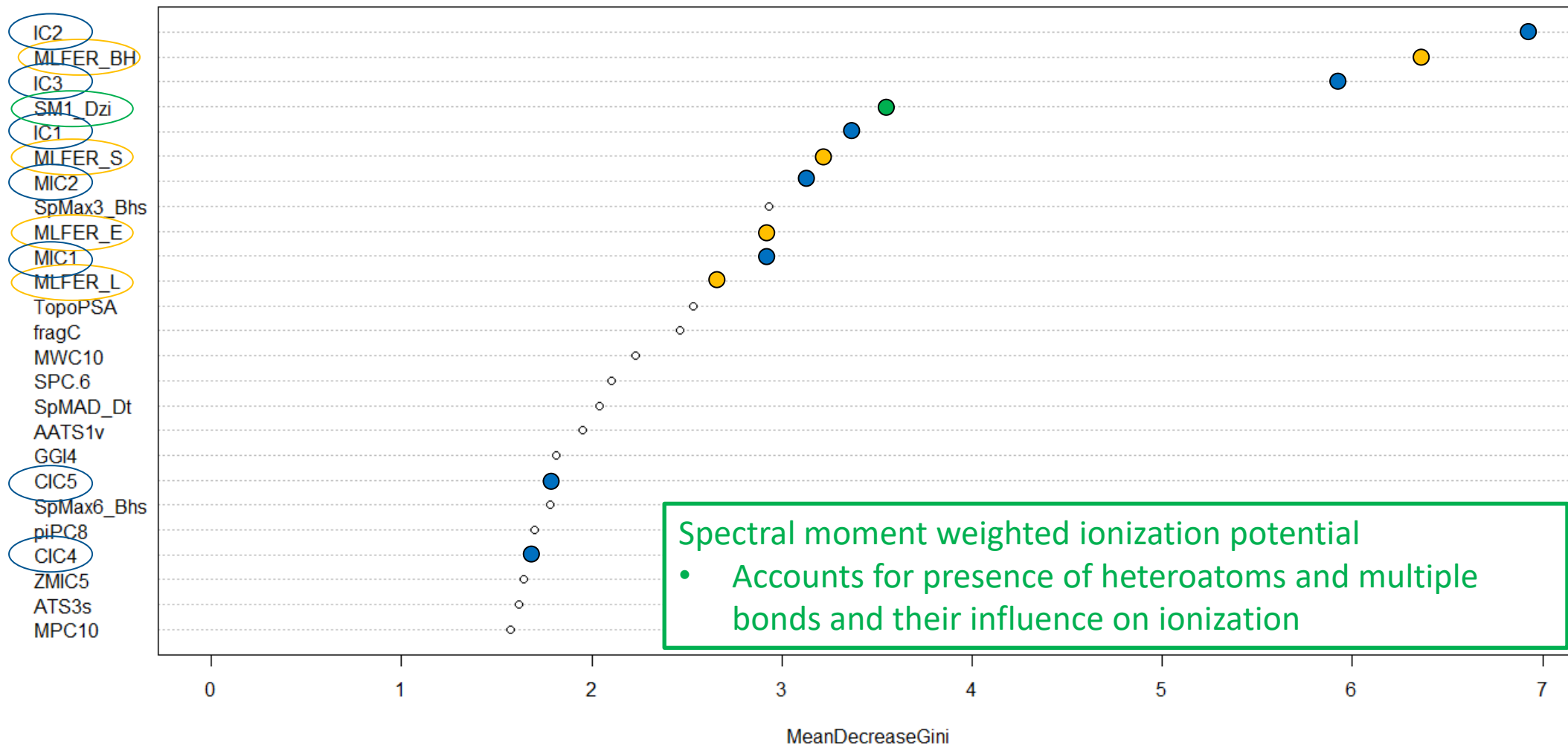
# Mechanistic interpretation

ESI+ Models (Downsampled Majority Class)



# Mechanistic interpretation

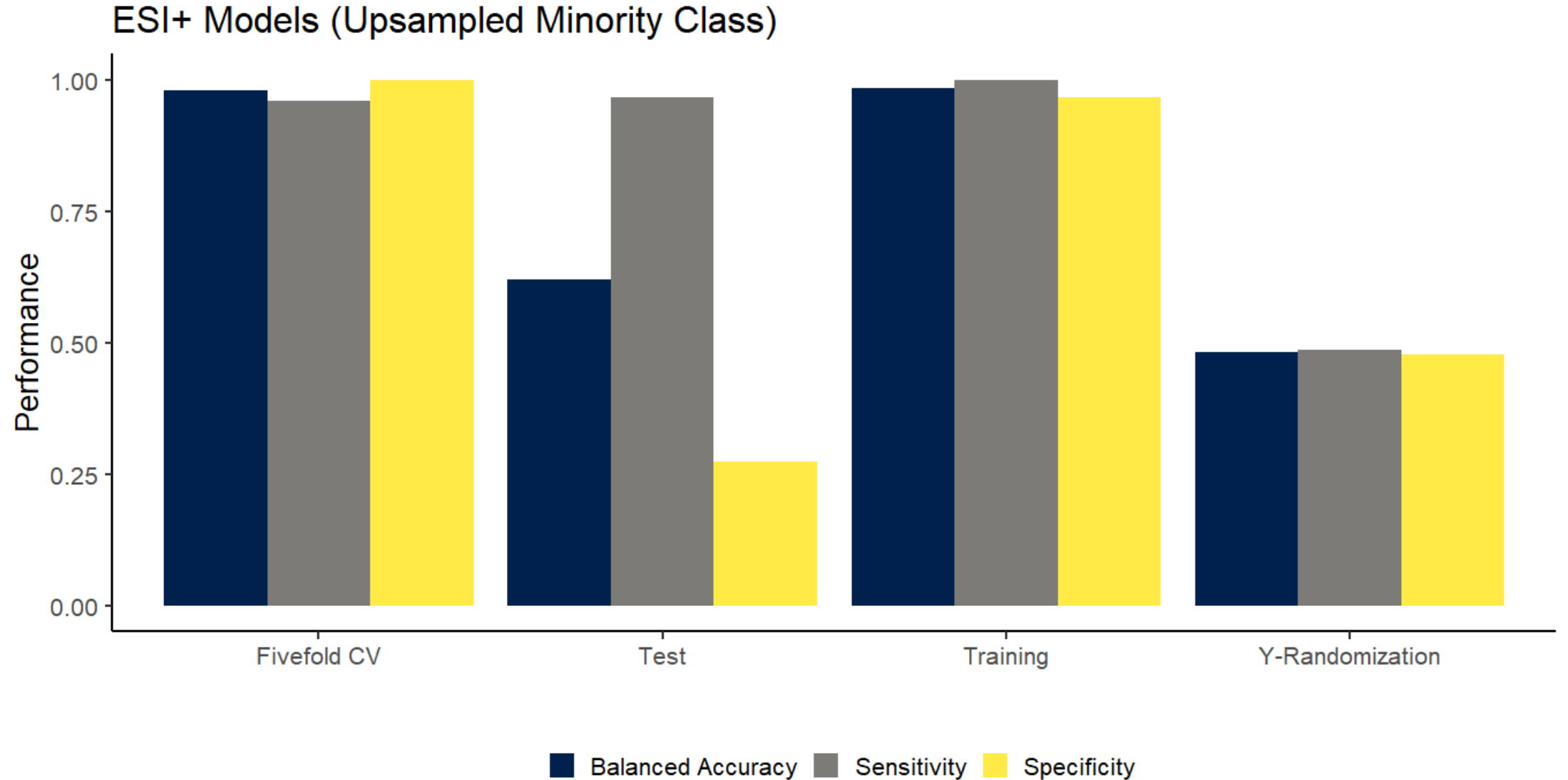
ESI+ Models (Downsampled Majority Class)



## Model summary

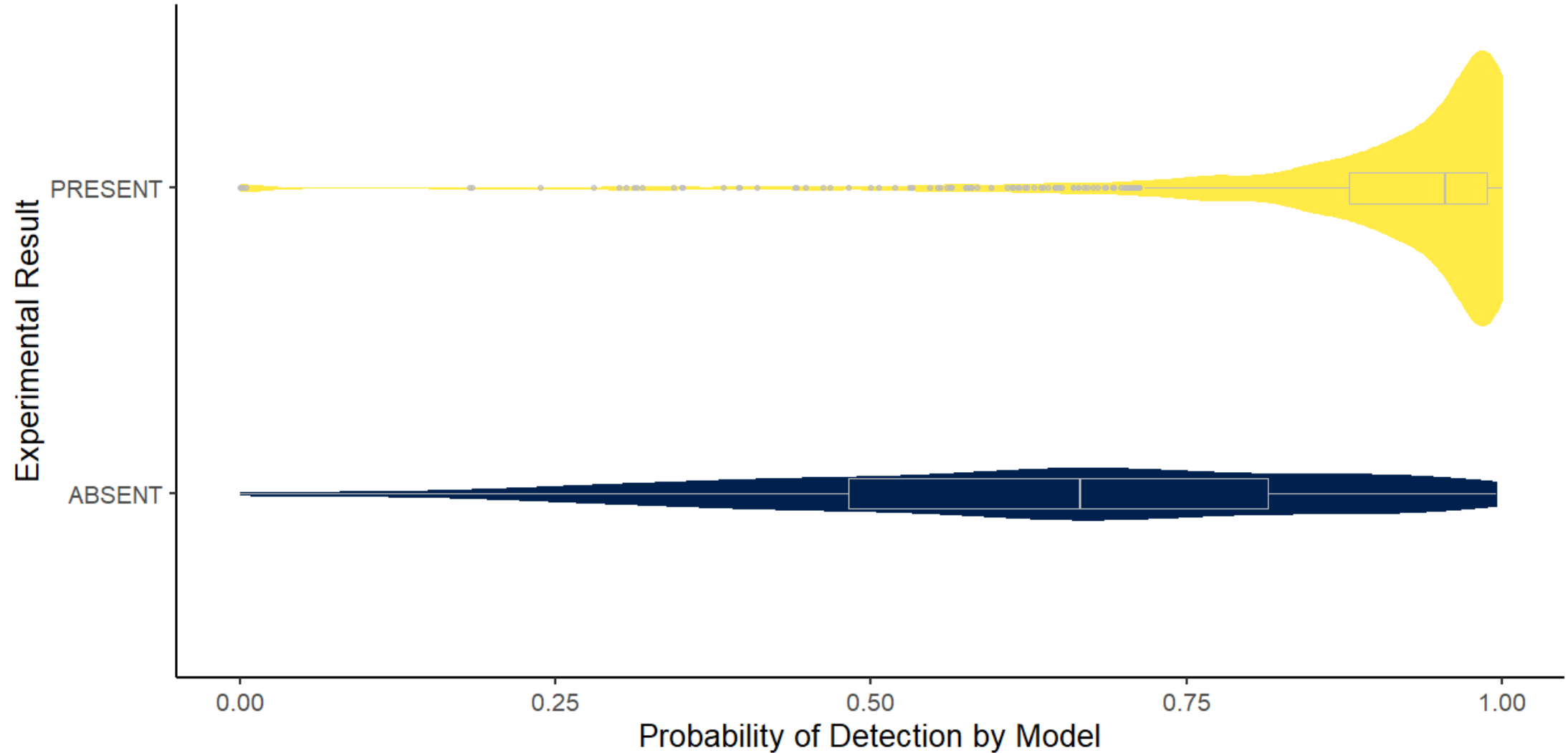
- Downsampled model provides excellent predictions for both amenable and unamenable compounds
  - Caveat: reduces sample space of amenable compounds
    - May not accurately predict every amenable compound
- Preferred model for ranking candidates in a suspect-screening analysis

# Model performance



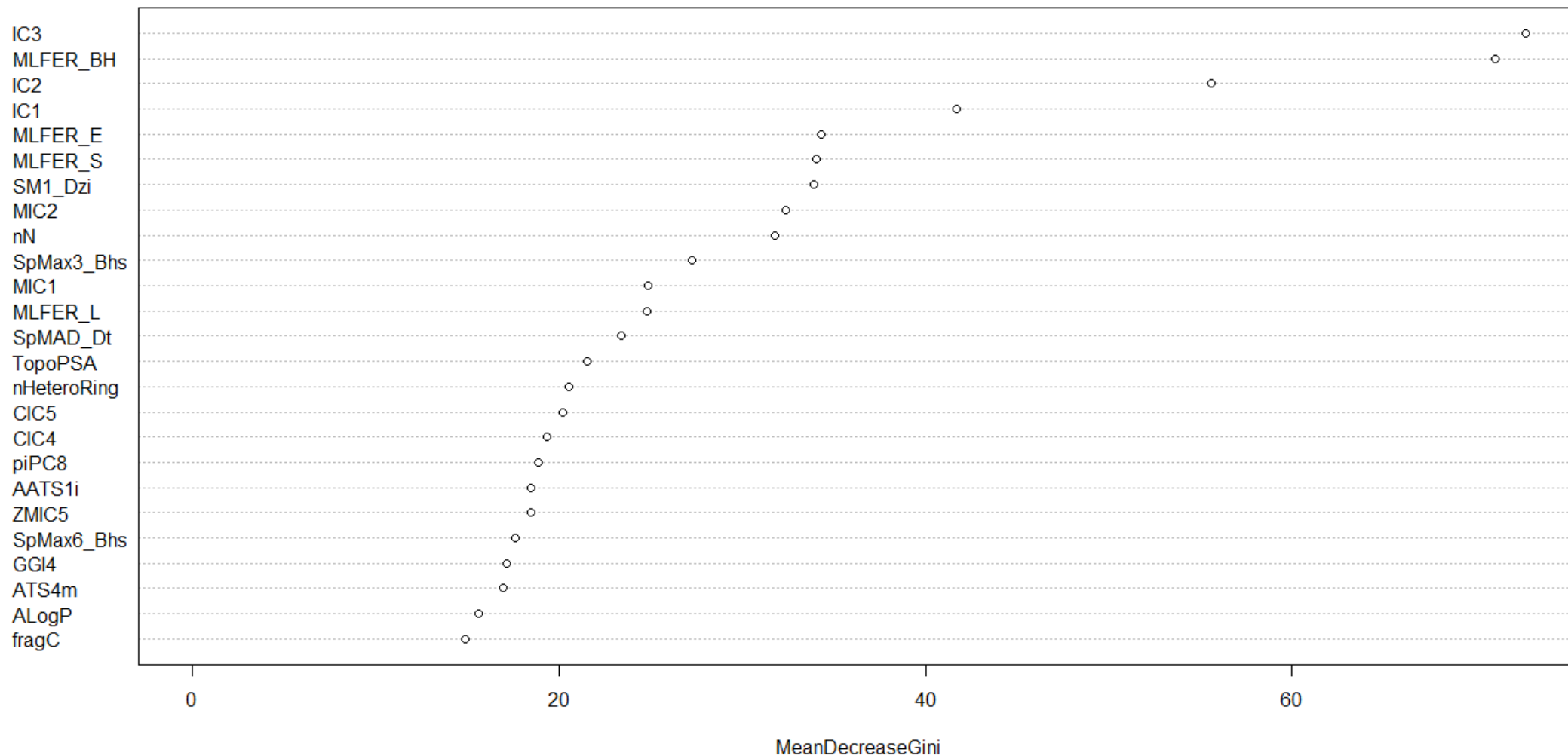
# Model performance

## ESI+ Models (Upsampled Minority Class)



# Model performance

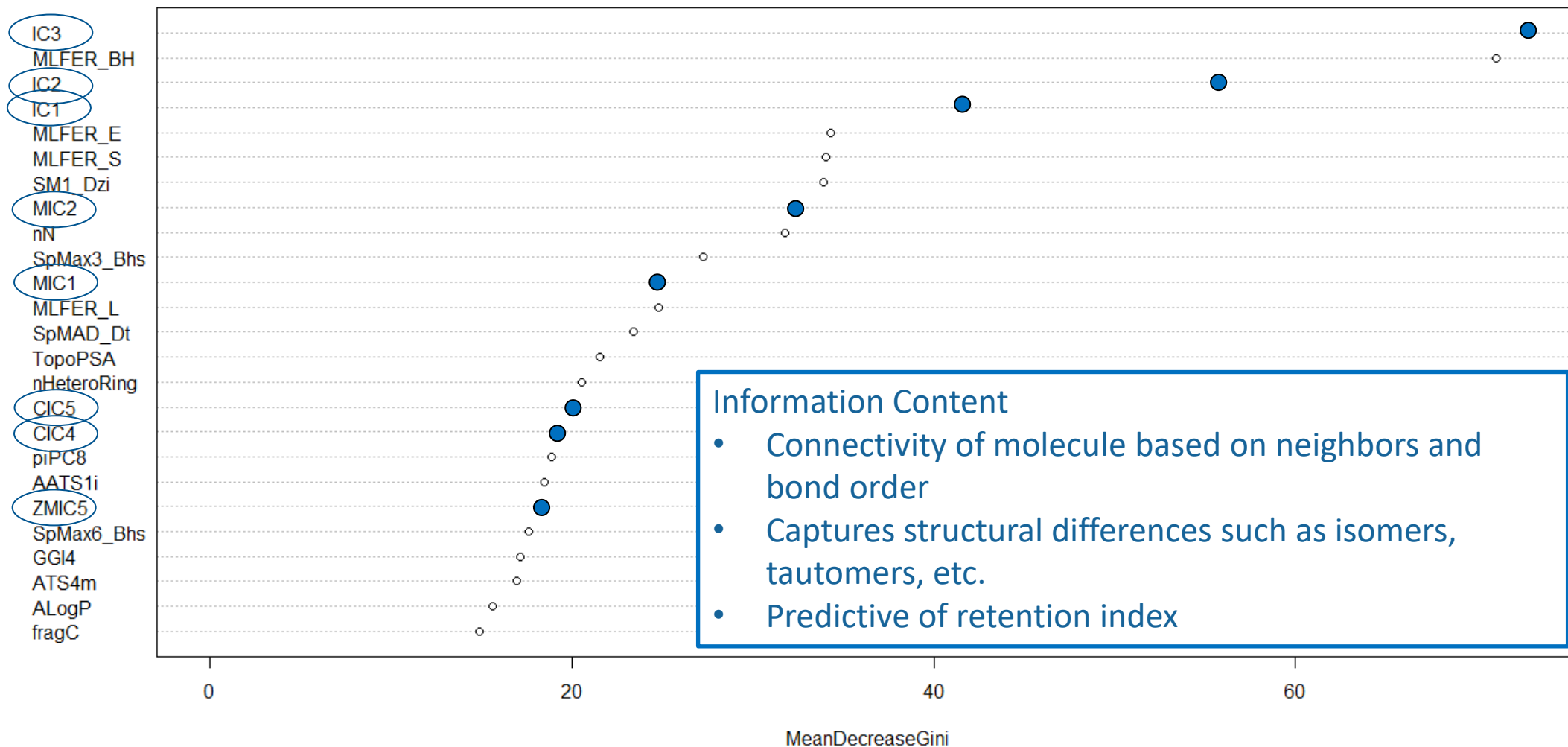
**ESI+ Models (Upsampled Minority Class)**





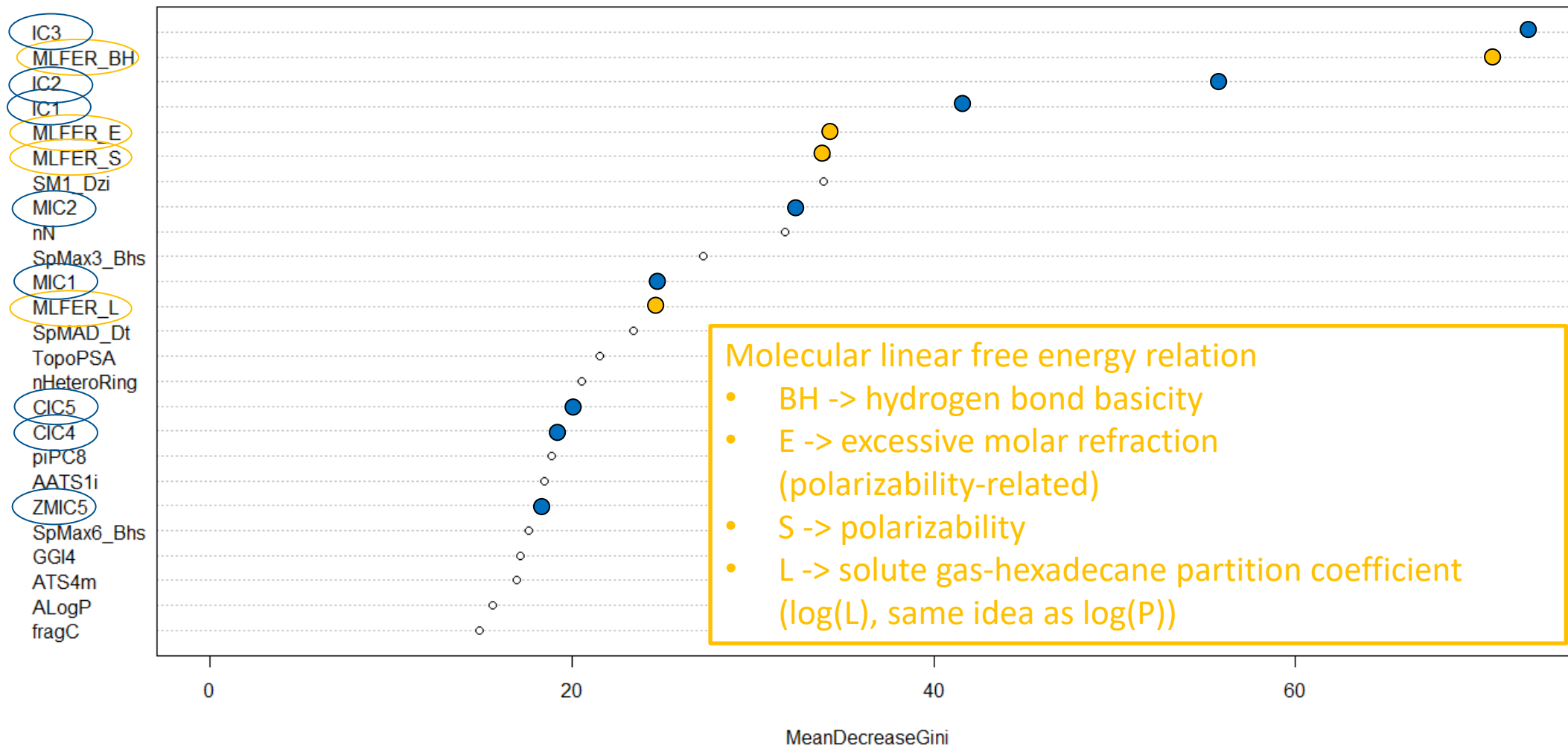
# Mechanistic interpretation

ESI+ Models (Upsampled Minority Class)



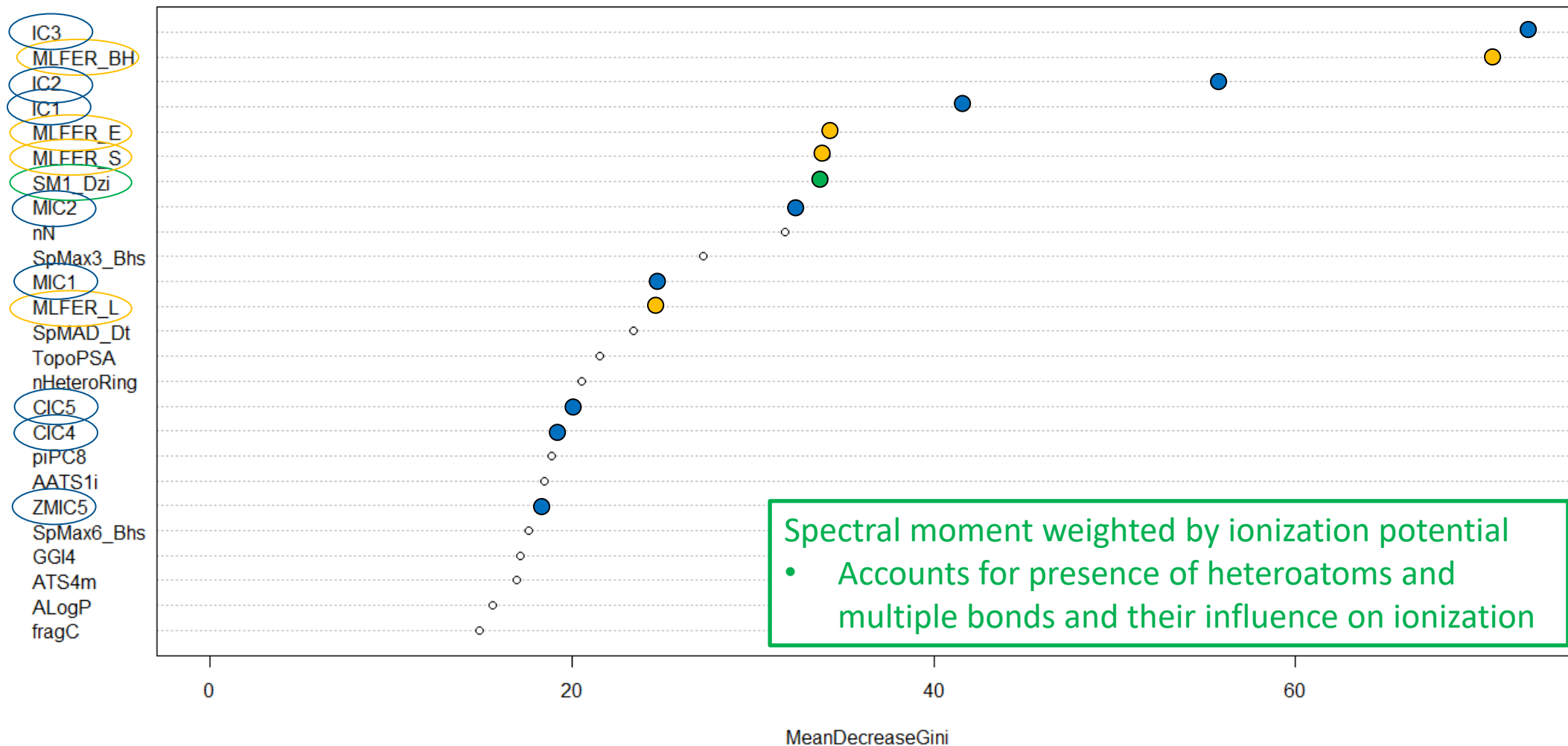
# Mechanistic interpretation

ESI+ Models (Upsampled Minority Class)



# Mechanistic interpretation

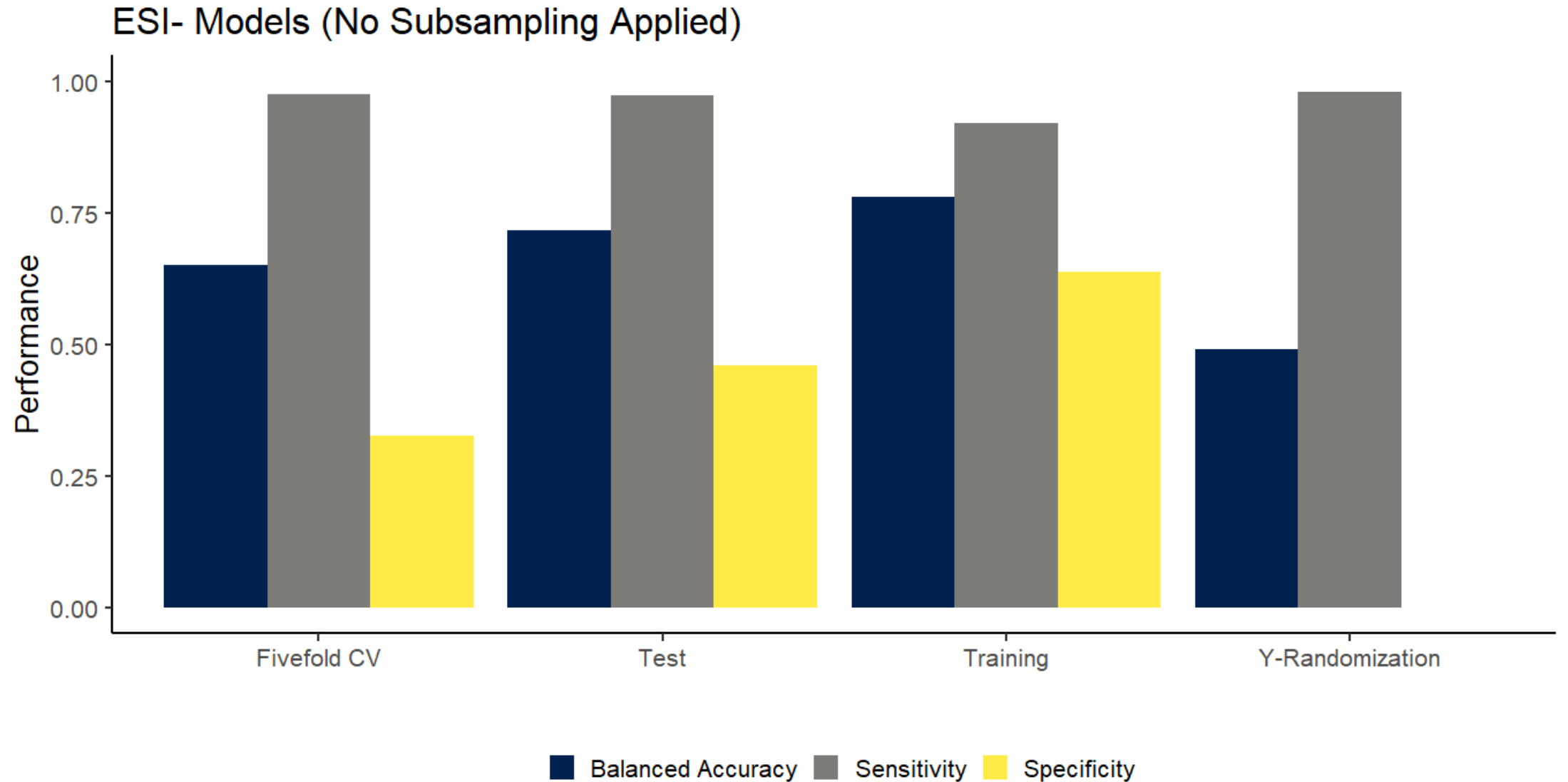
ESI+ Models (Upsampled Minority Class)



## Model summary

- Upsampled model provides excellent predictions for amenable compounds
  - Much larger sample space than downsampled model
  - Weak predictive power for unamenable compounds
  - Too optimistic for suspect-screening
- Preferred model for establishing which chemicals *may* be amenable to method
  - establishing a list of chemical standards

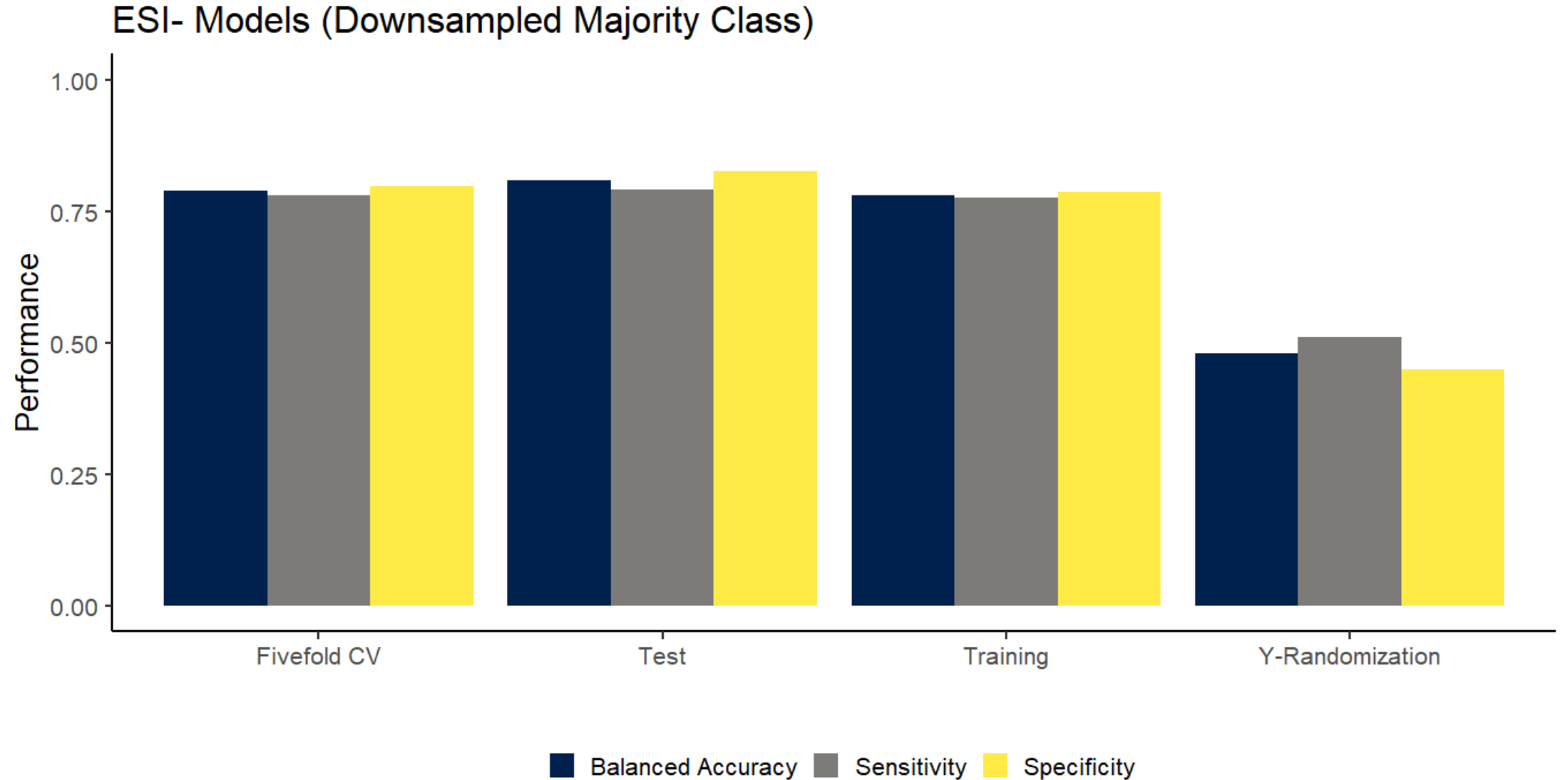
# Model performance



# Model summary

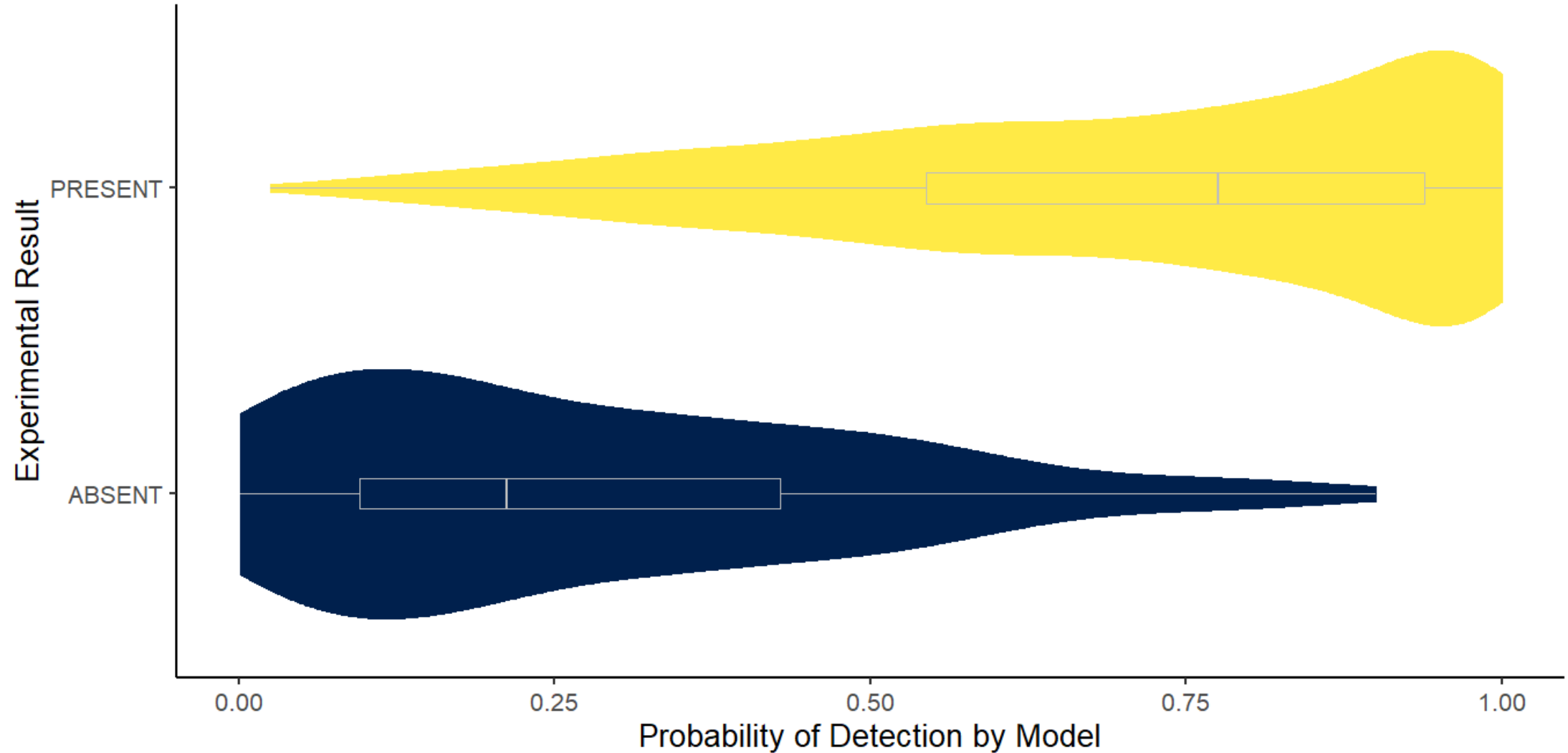
Imbalanced training sets lead to bad models!

# Model performance



# Model performance

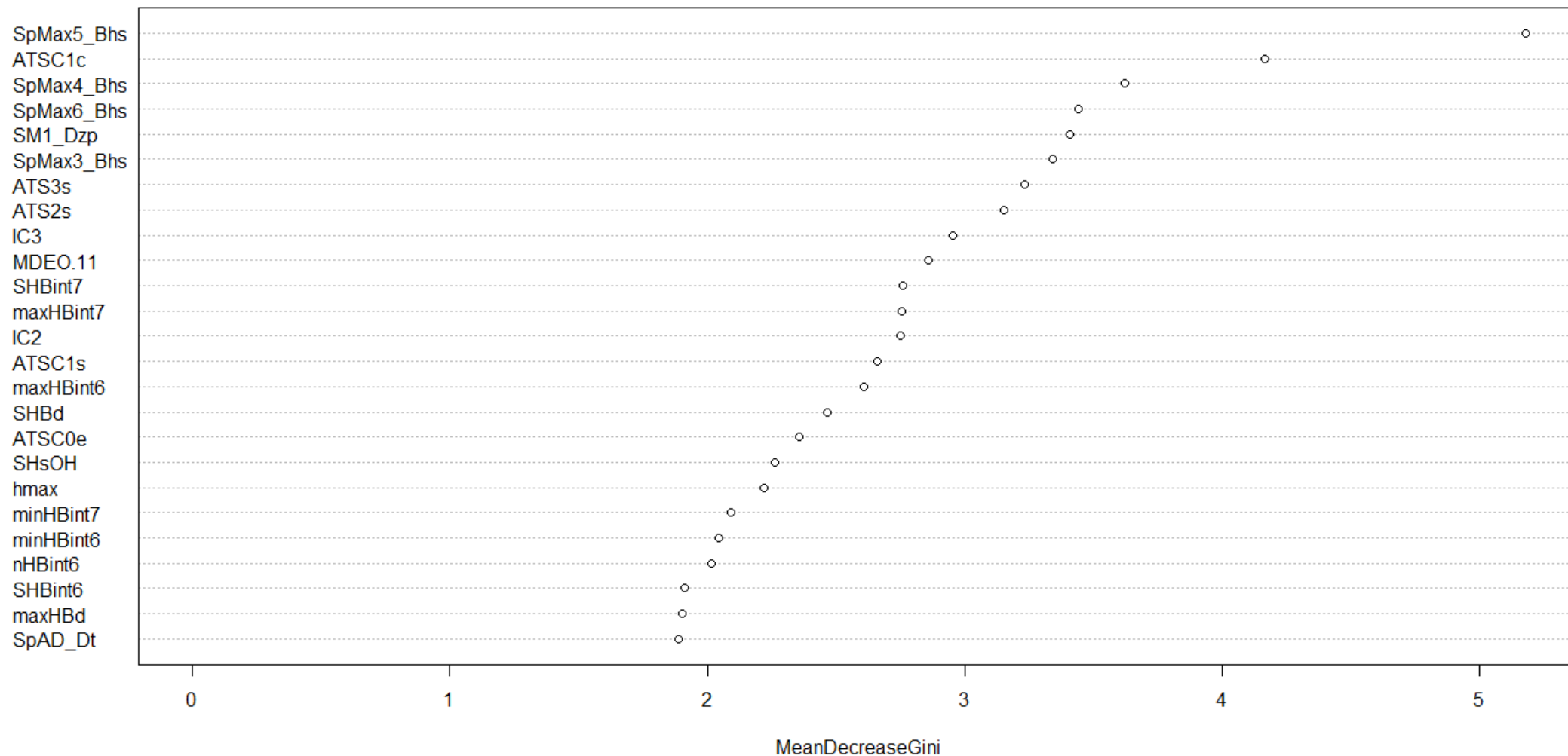
ESI- Models (Downsampled Majority Class)





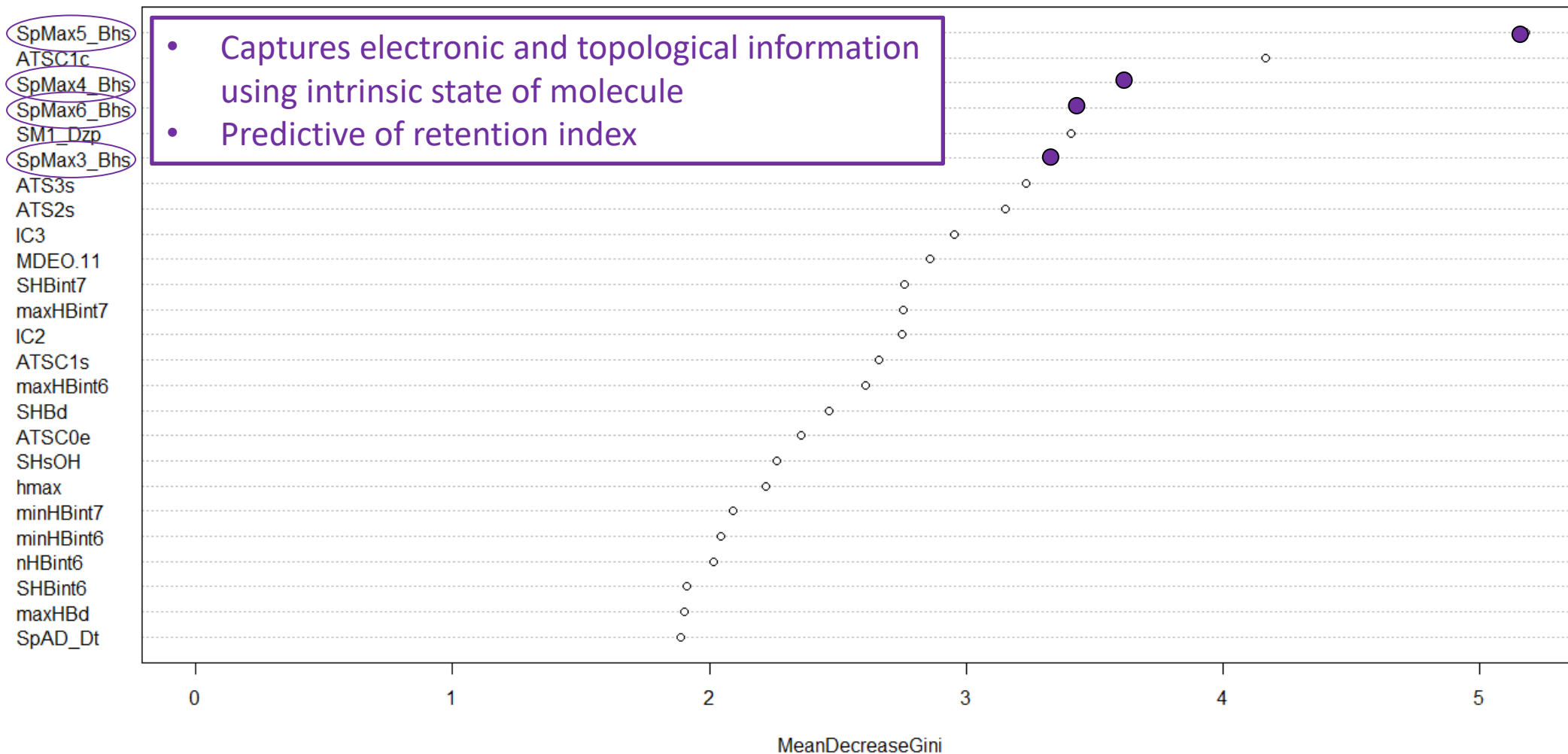
# Model performance

**ESI- Models (Downsampled Majority Class)**



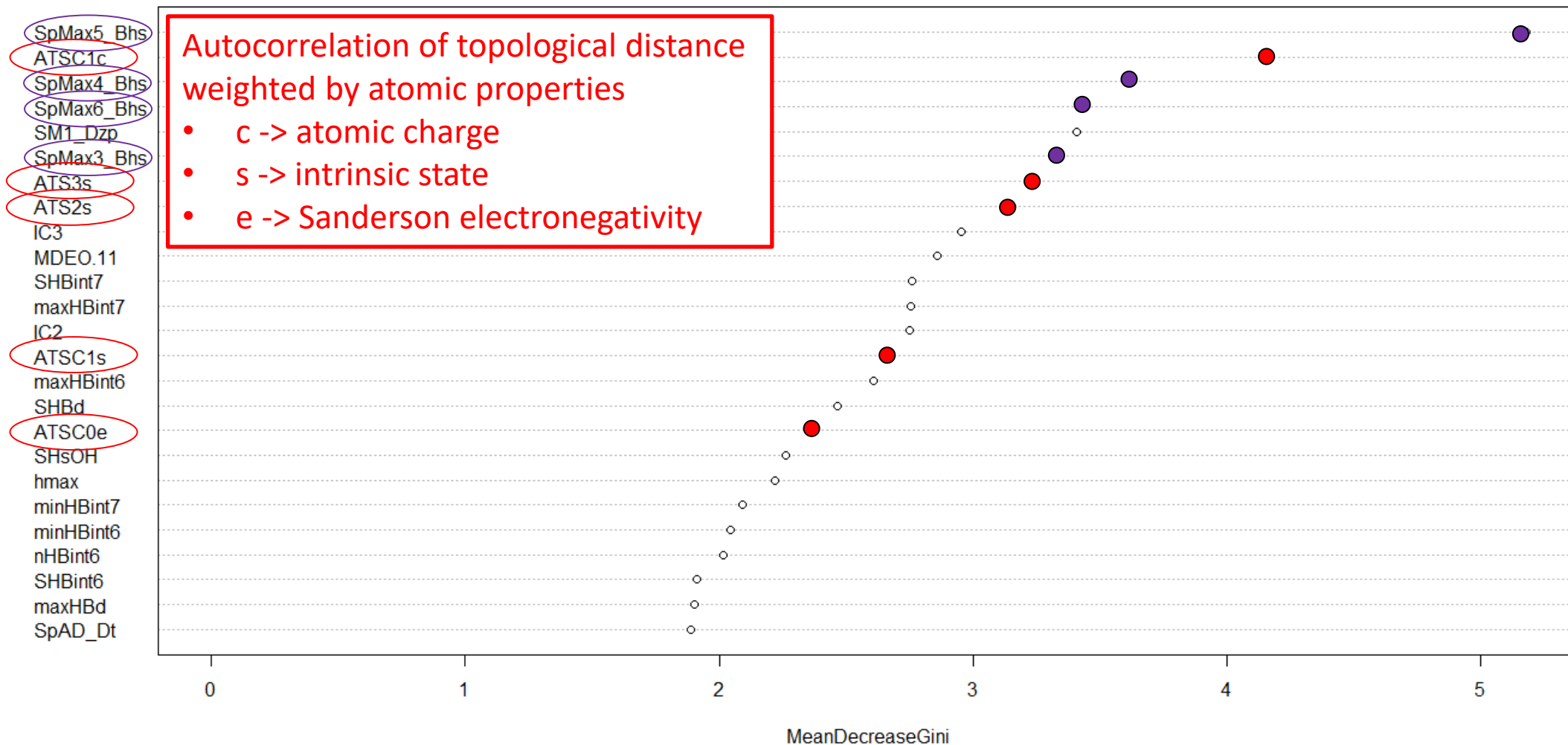
# Mechanistic interpretation

ESI- Models (Downsampled Majority Class)



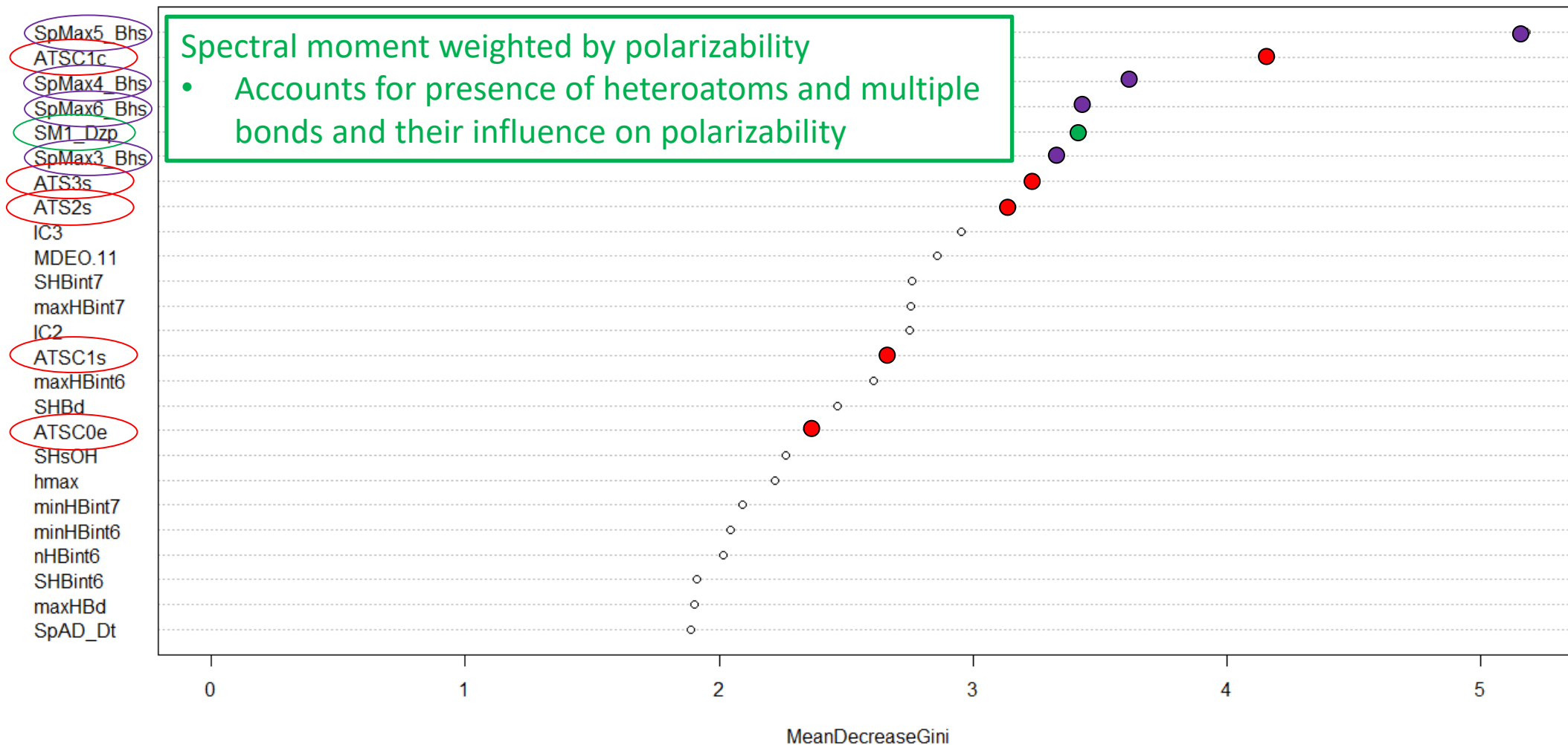
# Mechanistic interpretation

ESI- Models (Downsampled Majority Class)



# Mechanistic interpretation

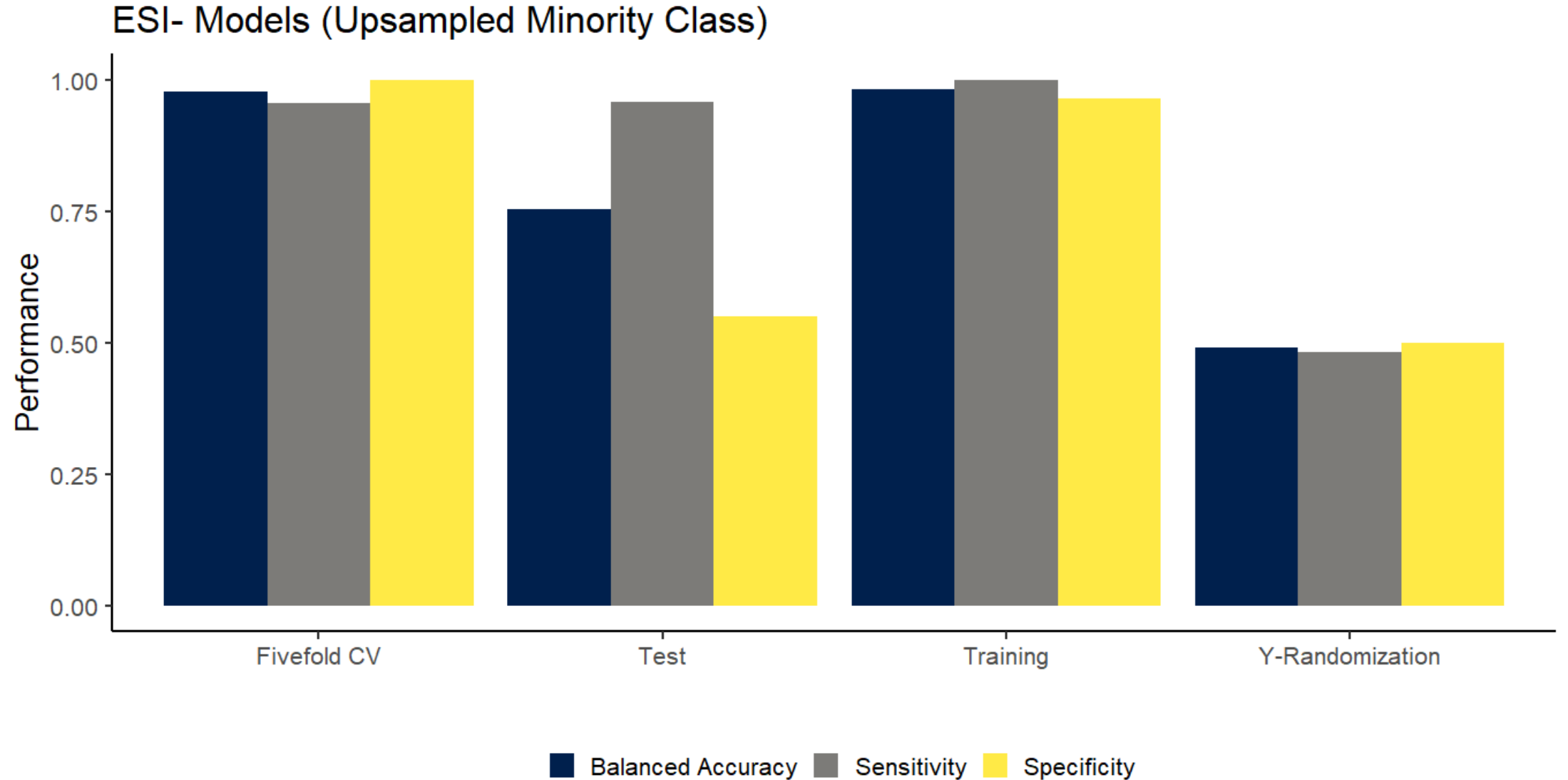
ESI- Models (Downsampled Majority Class)



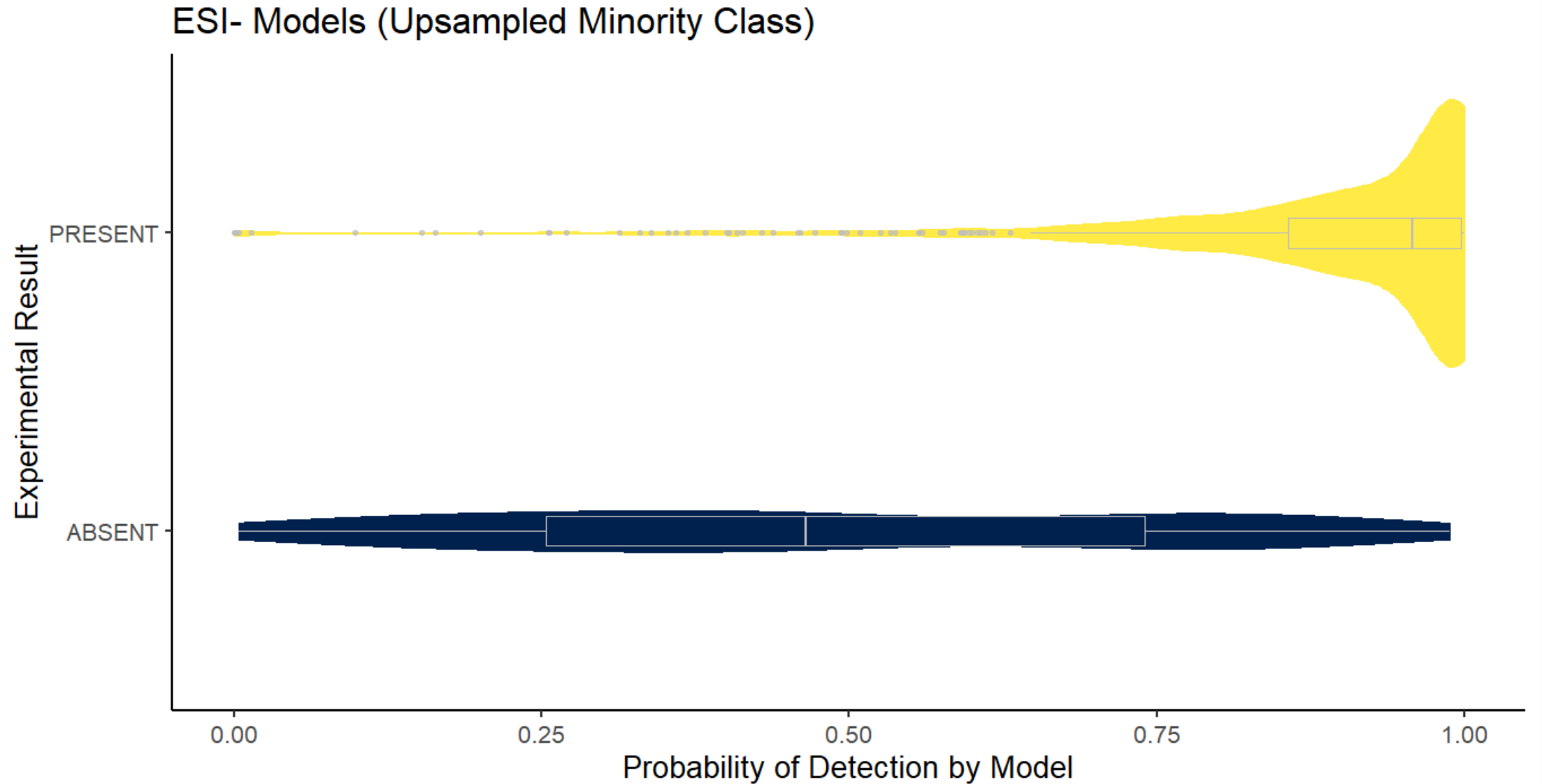
## Model summary

- Downsampled model provides excellent predictions for both amenable and unamenable compounds
  - Caveat: reduces sample space of amenable compounds
    - May not accurately predict every amenable compound
- Preferred model for ranking candidates in a suspect-screening analysis

# Model performance

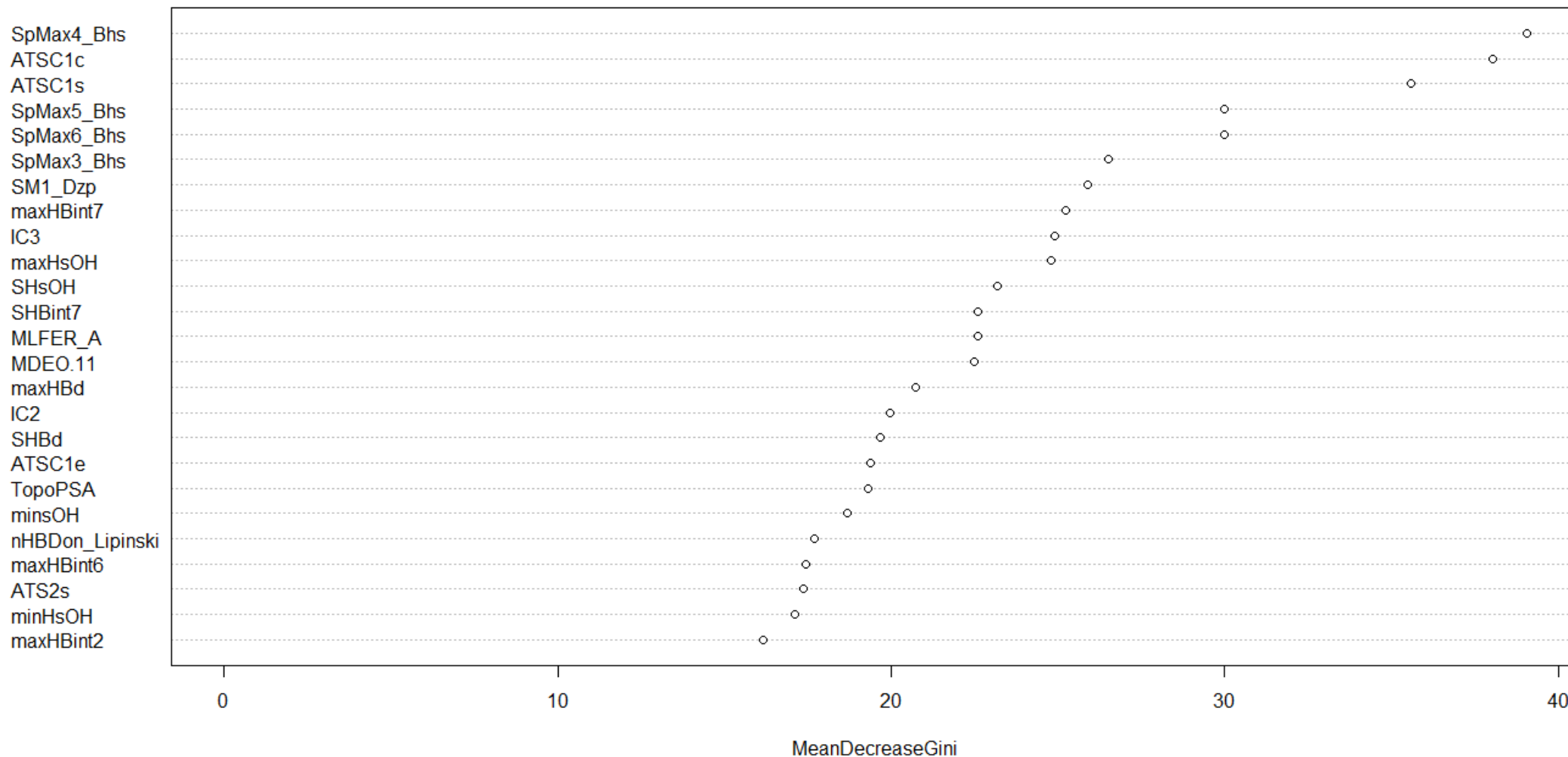


# Model performance



# Model performance

ESI- Models (Upsampled Minority Class)





# Mechanistic interpretation

ESI- Models (Upsampled Minority Class)

SpMax4\_Bhs

ATSC1c

ATSC1s

SpMax5\_Bhs

SpMax6\_Bhs

SpMax3\_Bhs

SM1\_Dzp

maxHBint7

IC3

maxHsOH

SHsOH

SHBint7

MLFER\_A

MDEO.11

maxHBd

IC2

SHBd

ATSC1e

TopoPSA

minsOH

nHBDOn\_Lipinski

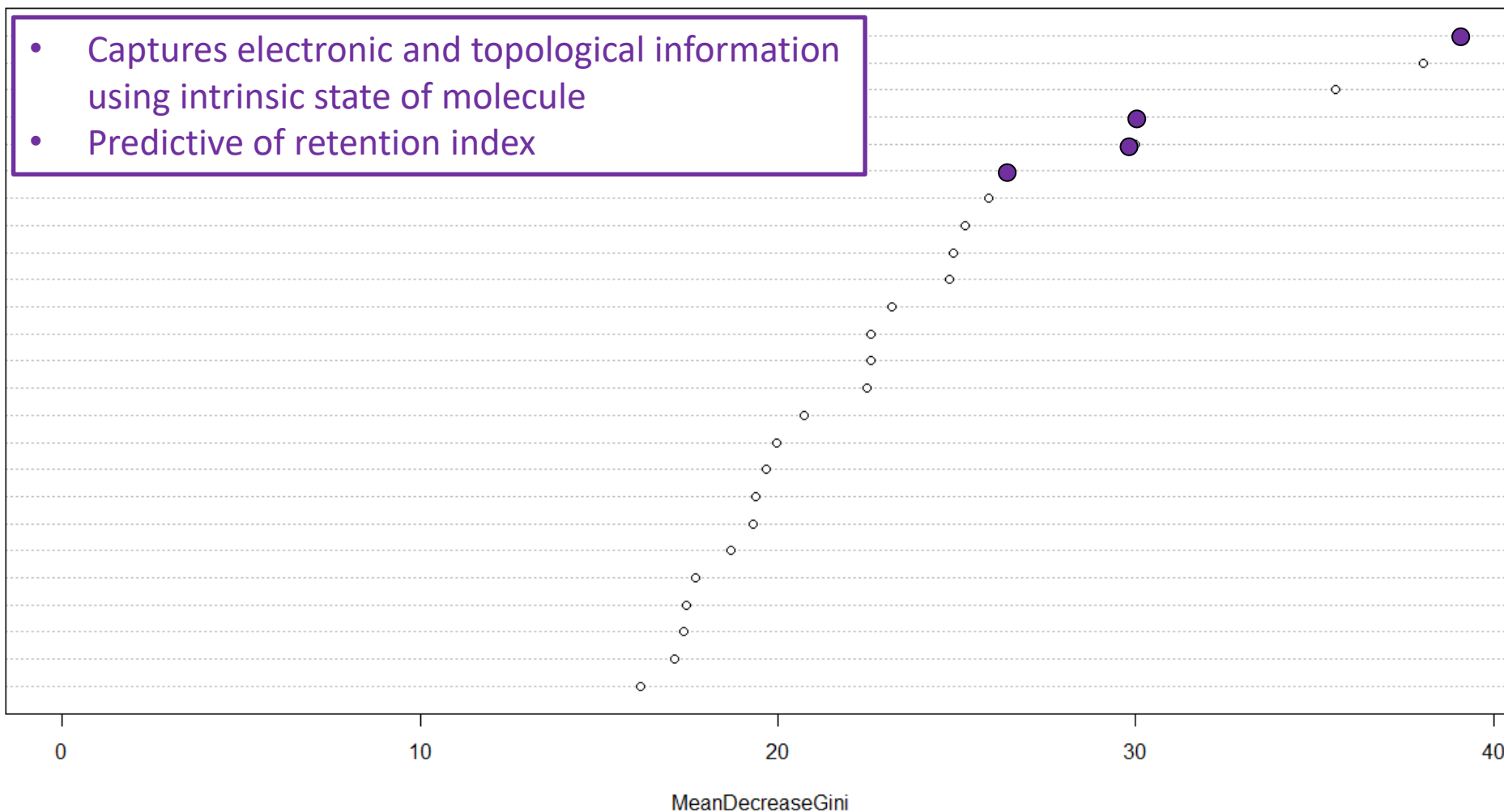
maxHBint6

ATS2s

minHsOH

maxHBint2

- Captures electronic and topological information using intrinsic state of molecule
- Predictive of retention index



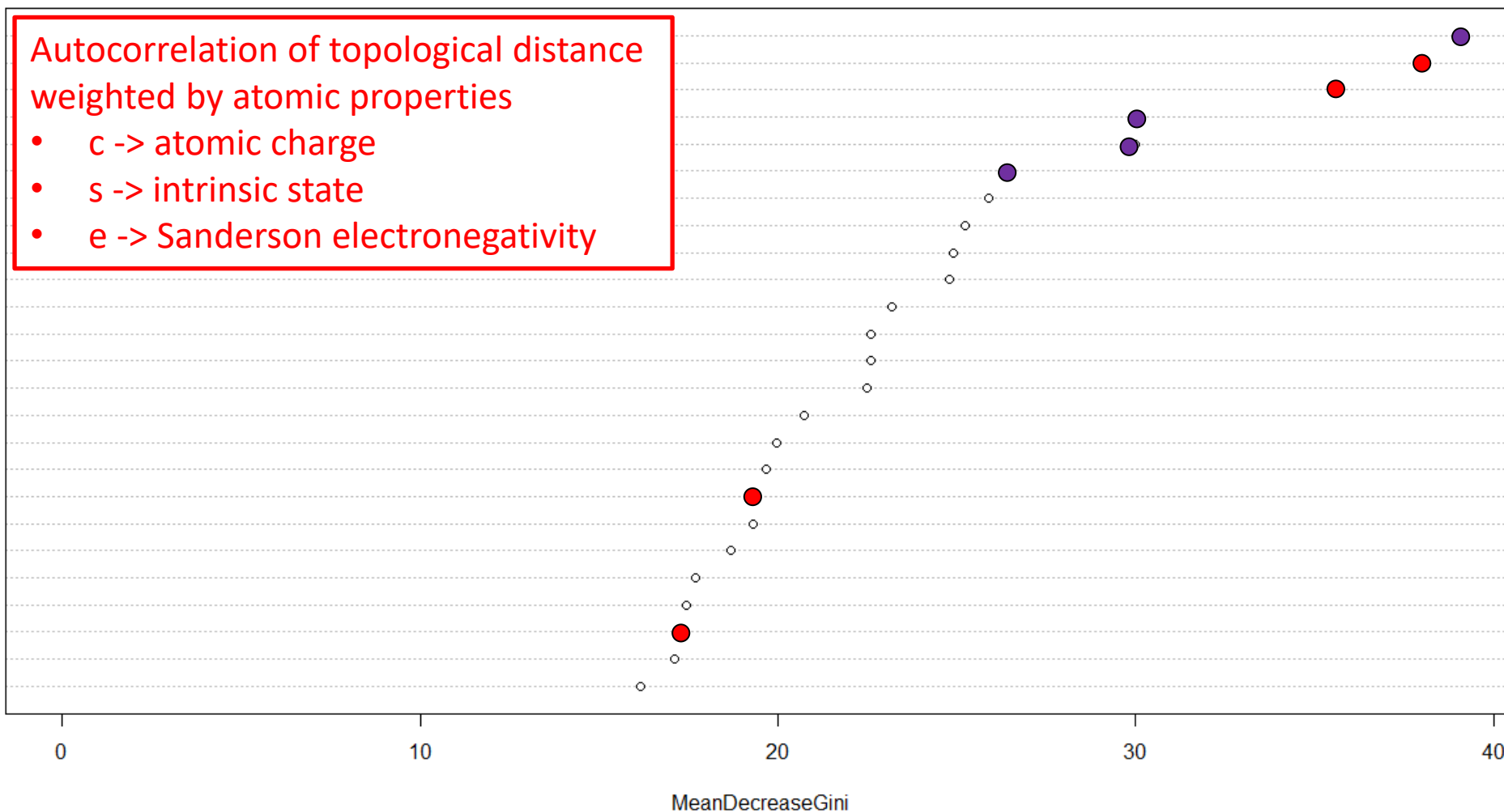
# Mechanistic interpretation

ESI- Models (Upsampled Minority Class)

SpMax4\_Bhs  
ATSC1c  
ATSC1s  
SpMax5\_Bhs  
SpMax6\_Bhs  
SpMax3\_Bhs  
SM1\_Dzp  
maxHBint7  
IC3  
maxHsOH  
SHsOH  
SHBint7  
MLFER\_A  
MDEO.11  
maxHBd  
IC2  
SHBd  
ATSC1e  
TopoPSA  
minsOH  
nHBDOn\_Lipinski  
maxHBint6  
ATS2s  
minHsOH  
maxHBint2

Autocorrelation of topological distance  
weighted by atomic properties

- c -> atomic charge
- s -> intrinsic state
- e -> Sanderson electronegativity



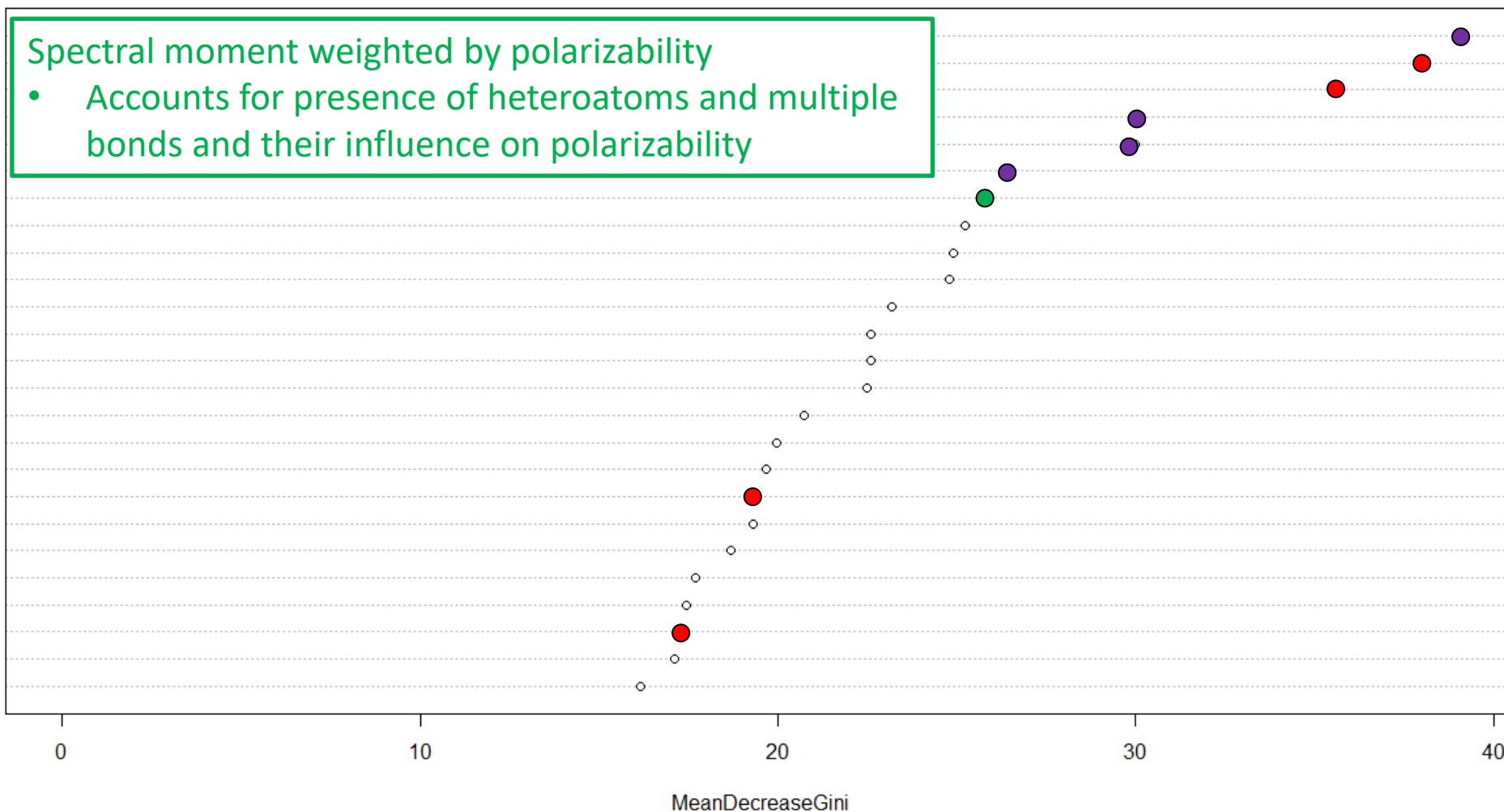
# Mechanistic interpretation

ESI- Models (Upsampled Minority Class)

SpMax4\_Bhs  
ATSC1c  
ATSC1s  
SpMax5\_Bhs  
SpMax6\_Bhs  
SpMax3\_Bhs  
SM1\_Dzp  
maxHBint7  
IC3  
maxHsOH  
SHsOH  
SHBint7  
MLFER\_A  
MDEO.11  
maxHBd  
IC2  
SHBd  
ATSC1e  
TopoPSA  
minsOH  
nHBDOn\_Lipinski  
maxHBint6  
ATS2s  
minHsOH  
maxHBint2

Spectral moment weighted by polarizability

- Accounts for presence of heteroatoms and multiple bonds and their influence on polarizability



## Model summary

- Upsampled model provides excellent predictions for amenable compounds
  - Much larger sample space than downsampled model
  - Weak predictive power for unamenable compounds
  - Too optimistic for suspect-screening
- Preferred model for establishing which chemicals *may* be amenable to method
  - establishing a list of chemical standards

## Current & future work

- Currently wrapping up manuscript
- Comparison of model results to Analytical QC data for ToxCast library
  - Good examples – no signal in LCMS ESI+, ESI- or in GCMS BUT present and high purity by NMR
- Compare model results to ENTACT results
  - Model predictions vs. independent labs, consensus of labs
- Future plans
  - Ensemble of upsampled and downsampled models?
  - Predictions for entirety DSSTox
  - Application for on-the-fly predictions based on a drawn structure

# CompTox Chemicals Dashboard mockup - Predictions

## Predictions

Chemical structure drawing tools (left sidebar)

Chemical structure drawing tools (bottom)

100% zoom

Calculate

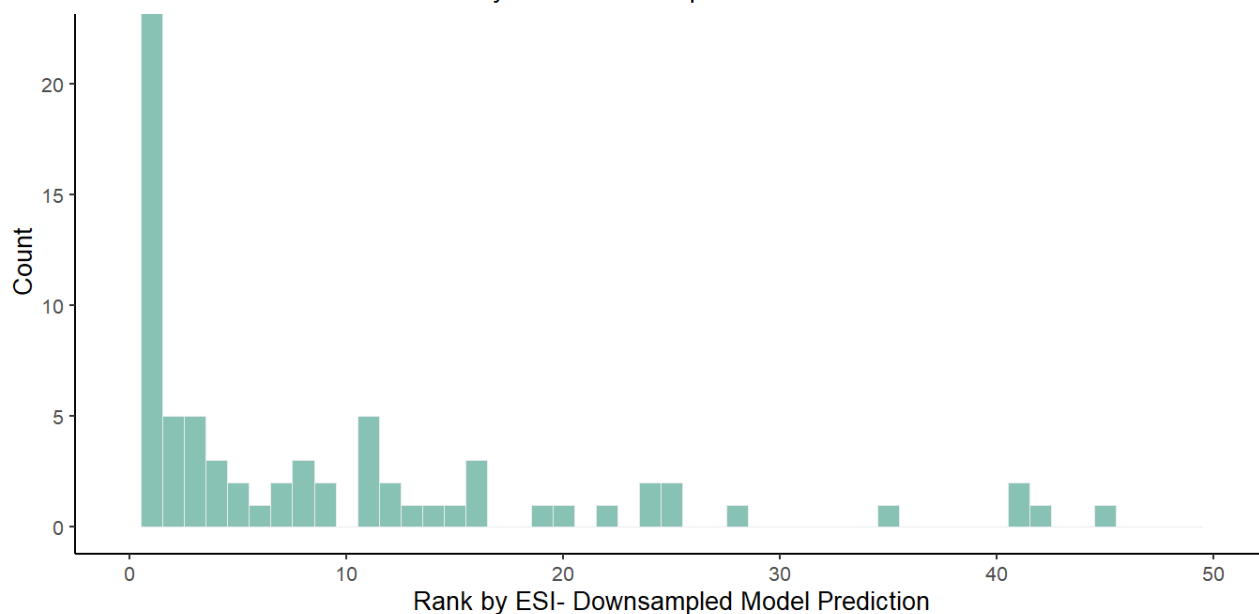
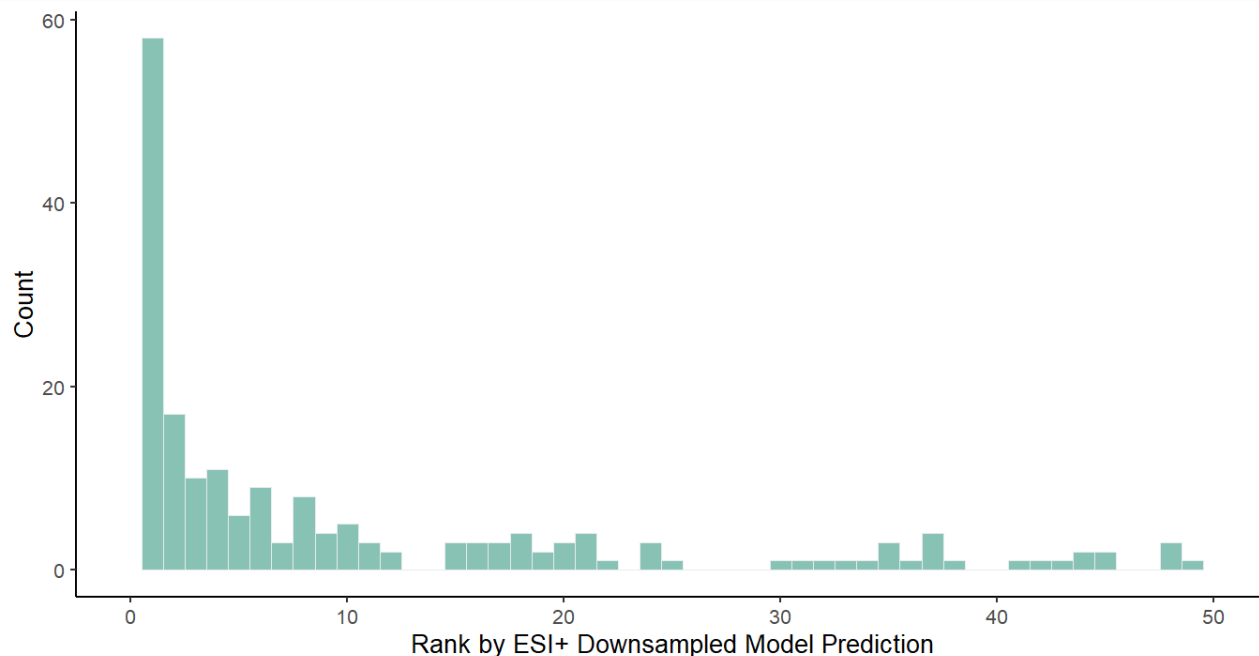
Select properties to predict

**T.E.S.T.**

- ☒ Toxicological properties
  - ☒ 96 hour fathead minnow LC50
  - ☒ 48 hour D. magna LC50
  - ☒ 48 hour T. pyriformis IGC50
  - ☒ Oral rat LD50
- ☒ Bioconcentration factor
- ☒ Developmental toxicity
- ☒ Ames mutagenicity
- ☒ Estrogen Receptor RBA
- ☒ Estrogen Receptor Binding
- ☒ Physical properties
  - ☒ Normal boiling point
  - ☒ Melting point
  - ☒ Flash point
  - ☒ Vapor pressure
  - ☒ Density
  - ☒ Surface tension
  - ☒ Thermal conductivity
  - ☒ Viscosity
  - ☒ Water solubility
- ☒ LCMS Predictions
  - ☒ ESI+ mode
  - ☒ ESI- mode

# Suspect-screening application

- List of ENTACT compounds identified in ESI+ & ESI- LCMS
  - 214 in ESI+
  - 105 in ESI-
- Retrieved candidates for each molecular formula via Dashboard
  - 13,325 candidates for ESI+
  - 7,079 candidates for ESI-
- Generated amenability predictions for candidate structures
- Rank ordered candidates by amenability probability





# Contributing researchers



## **EPA ORD**

Hussein Al-Ghoul\*  
Alex Chao\*  
Louis Groff\*  
Jarod Grossman\*  
Kristin Isaacs  
Sarah Laughlin\*  
Hannah Liberatore  
James McCord  
Kelsey Miller  
Jeff Minucci  
Seth Newton  
Katherine Phillips  
Allison Phillips\*  
Tom Purucker  
Randolph Singh\*  
Jon Sobus  
Mark Strynar  
Elin Ulrich  
Nelson Yeung\*

## **EPA ORD (cont.)**

Kathie Dionisio  
Chris Grulke  
Kamel Mansouri\*  
Andrew McEachran\*  
Ann Richard  
Adam Swank  
John Wambaugh  
Antony Williams

## **Agilent**

Jarod Grossman  
Andrew McEachran

## **GDIT**

Ilya Balabin  
Tom Transue  
Tommy Cathey

\* = ORISE/ORAU





Thank you for  
Listening!