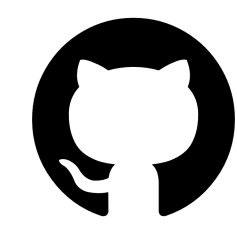


He He, Sheng Zha, Haohan Wang

<https://github.com/hhexiy/debiased>

## Overview

**Dataset bias:** spurious association between input and output

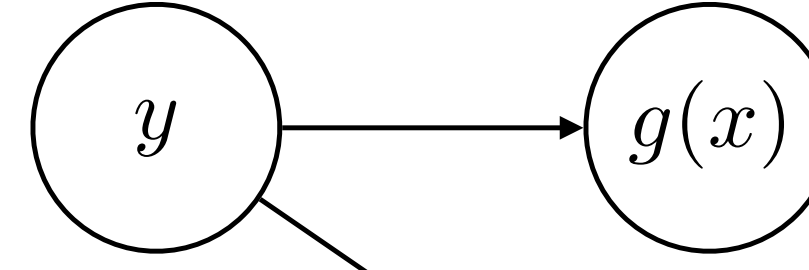
**Problem:** brittle model under slight distribution shift

**Goal:** Guard against *known* dataset bias

**Key idea:** don't learn from examples with strong (known) bias

### Training

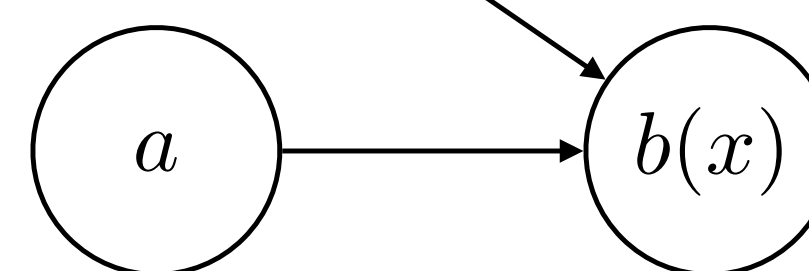
**Label:**  
contradiction



**Semantics:**  
P: The little girl is sad.  
H: The girl is not sad.



**Bias cause:**  
annotation strategy



**Word choice:**  
"not"

$p(y, b(x))$  may change at test time

## Debiasing by Residual Fitting (DRiFt)

1. Learn a **biased classifier** using **features** that should not be associated with the label

$$\theta^* = \arg \min_{\theta} \mathbb{E}_P[L(f_s(I(x); \theta), y)]$$

2. Learn the **debiased classifier** by fitting the residuals

$$\min_{\phi} \mathbb{E}_P[L(f_s(I(x); \theta^*) + f_d(x; \phi), y)]$$

Fixed

Learns what cannot be explained by  $I(x)$

### Cross-entropy loss

Biased classifier:  $p_s(y | x) \propto \exp(f_s^y)$

Debiased classifier:  $p_d(y | x) \propto \exp(f_d^y)$

Learned by MLE using  $I(x)$

Objective:

$$J(\phi) = \sum_{(x,y)} \log p(y | x)$$

$$p(y | x) \propto \exp(f_s^y + f_d^y) \propto p_s p_d$$

$J_{MLE}$

Regularizer  $R(x)$

$$= C + \sum_{(x,y)} \left[ \log p_d(y | x) - \log \sum_k p_s^*(k | x) p_d(k | x) \right]$$

Predictive biased classifier:

$$p_s^* \rightarrow 1 \Rightarrow \nabla_{\phi} R(x) = -\nabla_{\phi} \log p_d(y | x)$$

Cancels MLE gradient

Uninformative biased classifier:

$$p_s = 1/K \Rightarrow p(y | x) = p_d(y | x)$$

Reduced to MLE

## Experiments

### Biased models:

- Hypothesis-only
- Handcrafted
- CBOW

### Debiased models:

- Finetuned BERT [Devlin+ 18]
- DA [Parikh+ 16]: ~BoW
- ESIM [Chen+ 17]: ~DA + LSTM

### Learning algorithms:

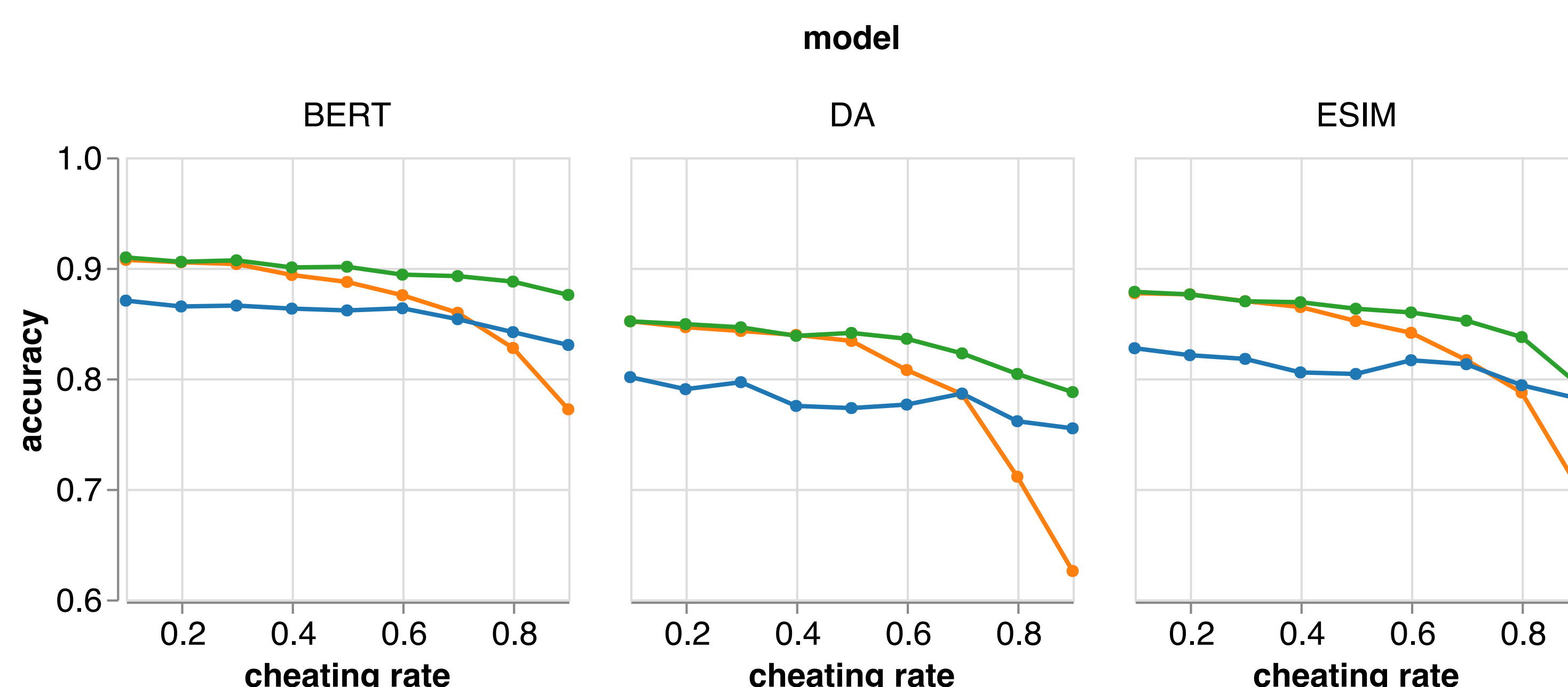
- MLE
- DRiFt
- Rm

Remove biased examples

### Synthetic bias:

Groundtruth with  $p_{\text{cheat}}$  at train and random at test  
P: I love dogs.  
H: [contradiction] I don't love dogs.

### Results on synthetic bias



MLE: baseline

DRiFt-hypo: hypothesis-only biased classifier

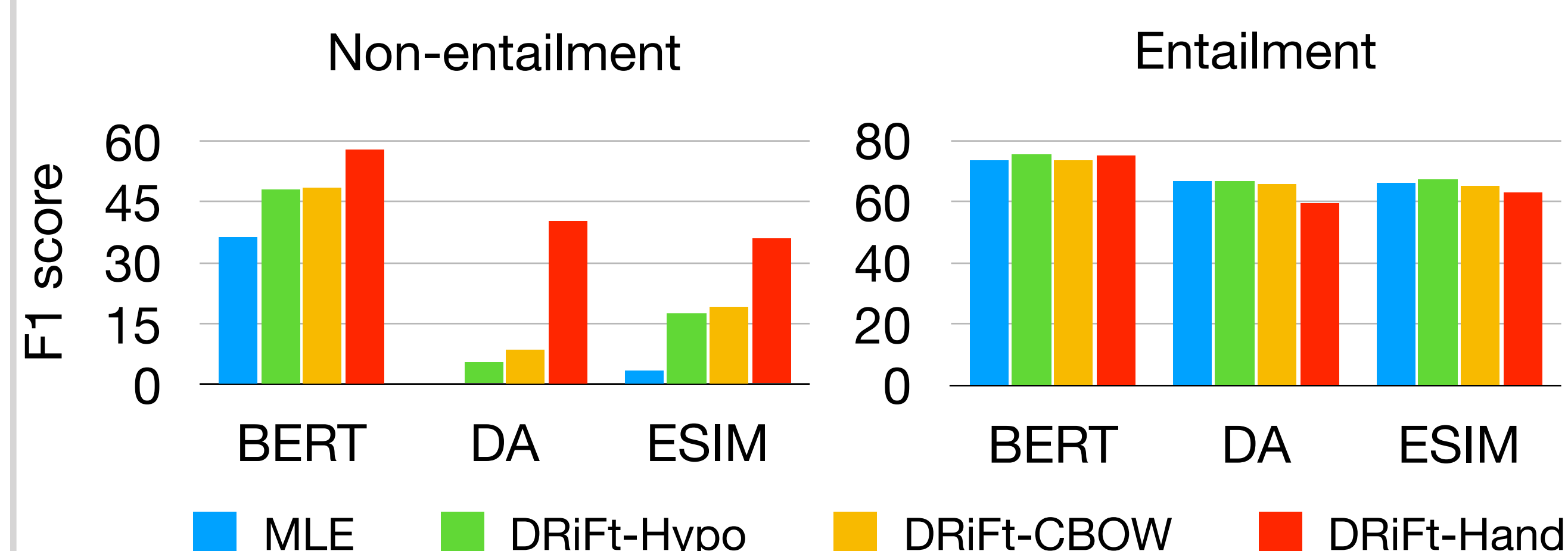
Rm-cheat: remove cheatable examples (oracle)

- Bias needs to be presented on a majority of examples
- BERT is more robust than non-pretrained models
- Importance of accurate prior knowledge on biased features (compare DRiFt-hypo and Rm-cheat)

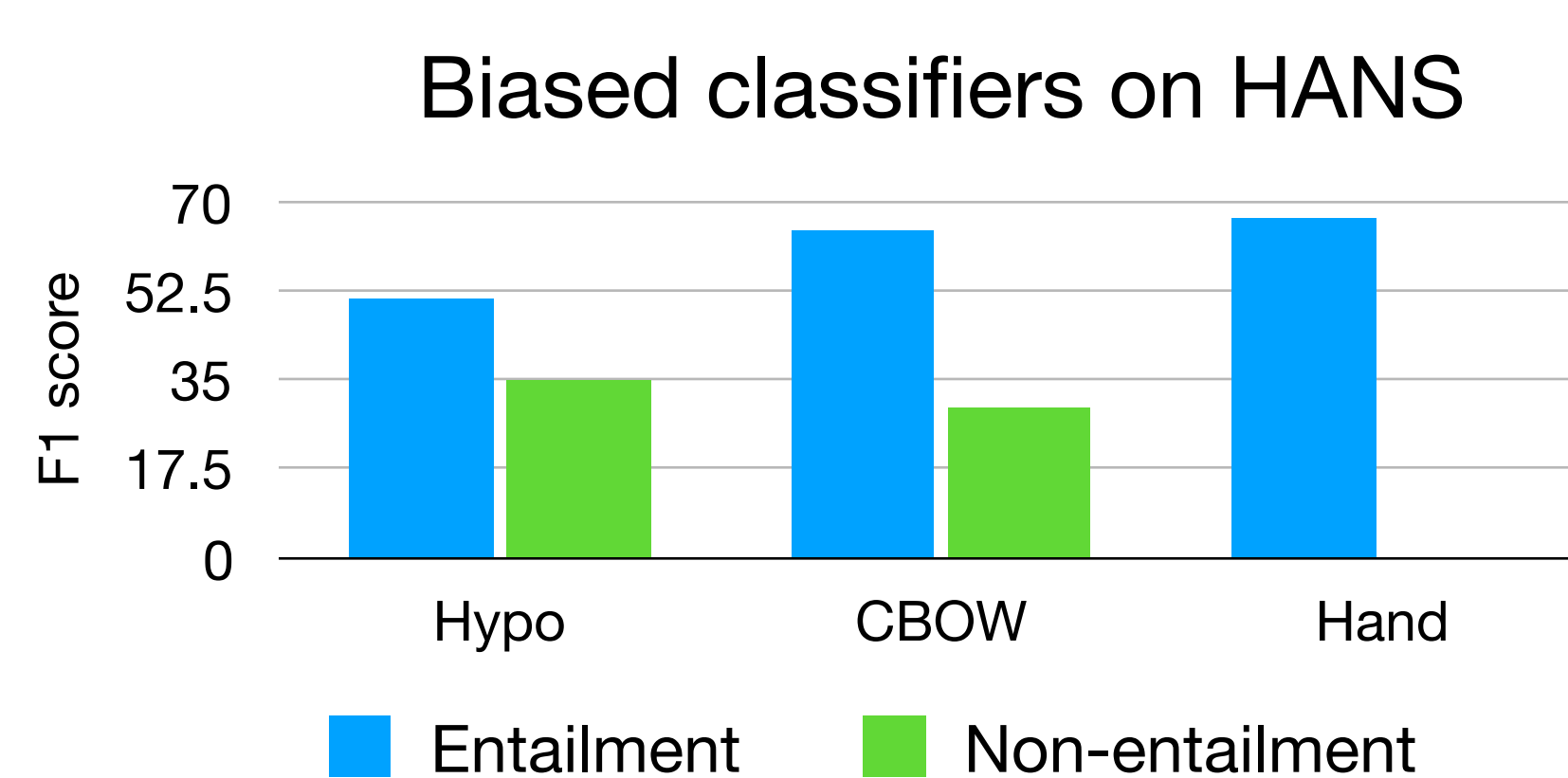
### NLI Results

**HANS** [McCoy+ 19]: exploit word overlap bias, e.g., The doctor was paid by the actor  $\Rightarrow$  The doctor paid the actor

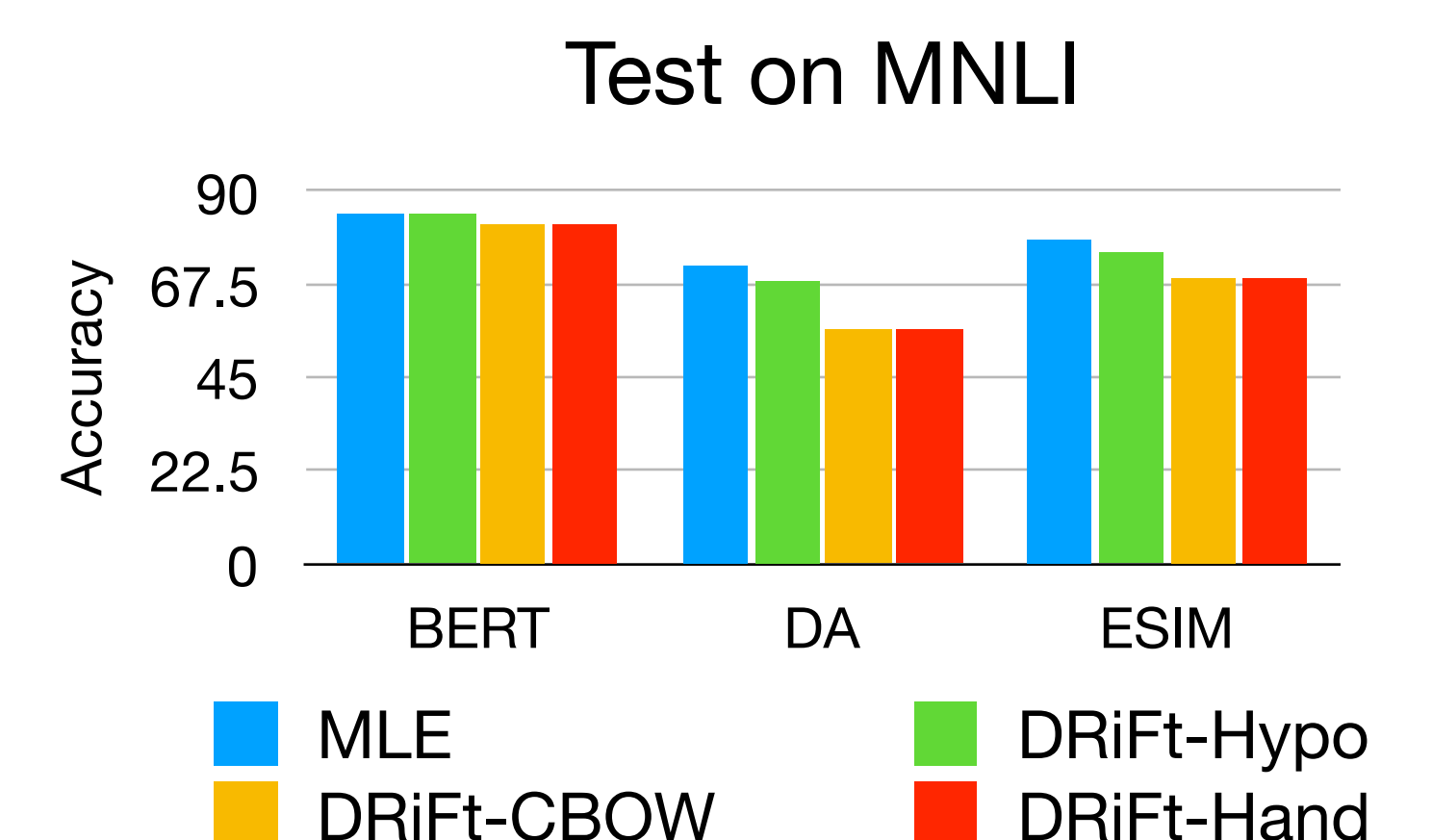
### Train on MNLI and test on HANS



### Why is "Hand" more effective?



### In-distribution accuracy



DRiFt improves performance on challenge data (non-entailment)

"Hand" better captures bias exploited by HANS

Accuracy-robustness trade-off