



THE OHIO STATE  
UNIVERSITY

# Identifying inherent disagreement in natural language inference

Xinliang Frederick Zhang and Marie-Catherine de Marneffe  
The Ohio State University

NAACL 2021

# Natural Language Inference (NLI)

**Premise:** A homeless man being observed  
by a man in business attire.

**Hypothesis:** Two men are sleeping in a hotel.



Contradiction



Neutral



Entailment

# Natural Language Inference (NLI)

**Premise:** A homeless man being observed by a man in business attire.

**Hypothesis:** Two men are sleeping in a hotel.



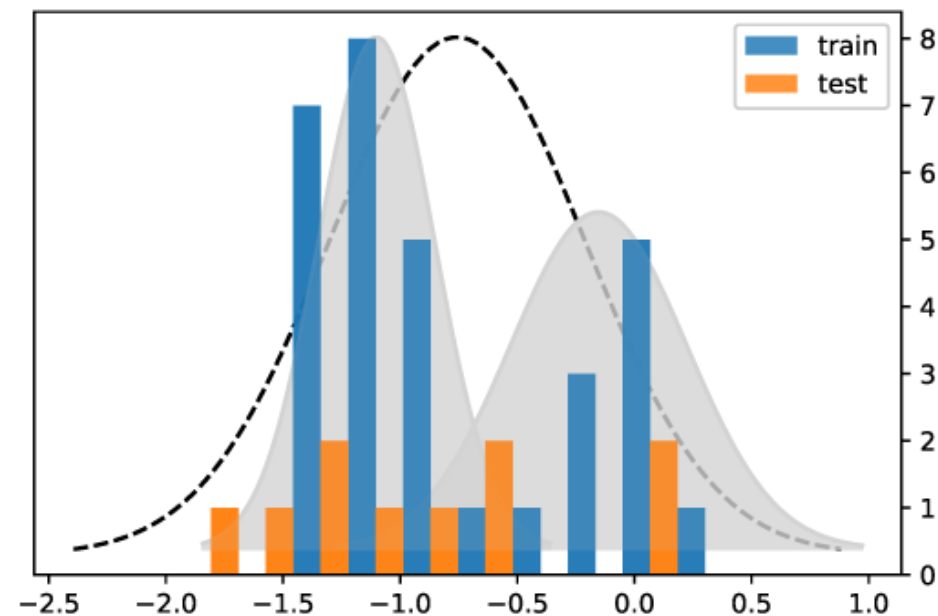
Contradiction



Neutral



Entailment



Contradiction Neutral





# Data: CommitmentBank

**Premise:** B: Yeah, and EDS is very particular about this, hair cuts, A: Wow. B: I mean it was like you can't have, you know, such and such facial hair, no beards, you know, and just really detailed. A: I don't know that **that would be a good environment to work in.**

**Hypothesis:** that would be a good environment to work in

**Label?** [ 2, 0, 0, 0, 0, -1, -2, -3 ]



# Data: CommitmentBank

**Premise:** B: Yeah, and EDS is very particular about this, hair cuts, A: Wow. B: I mean it was like you can't have, you know, such and such facial hair, no beards, you know, and just really detailed. A: I don't know that **that would be a good environment to work in.**

**Hypothesis:** that would be a good environment to work in

Label? [ 2, 0, 0, 0, 0, -1, -2, -3 ]





# Finer-grained labels for NLI

**Premise:** B: Yeah, and EDS is very particular about this, hair cuts, A: Wow. B: I mean it was like you can't have, you know, such and such facial hair, no beards, you know, and just really detailed. A: I don't know that **that would be a good environment to work in.**



Entailment



Neutral



Contradiction

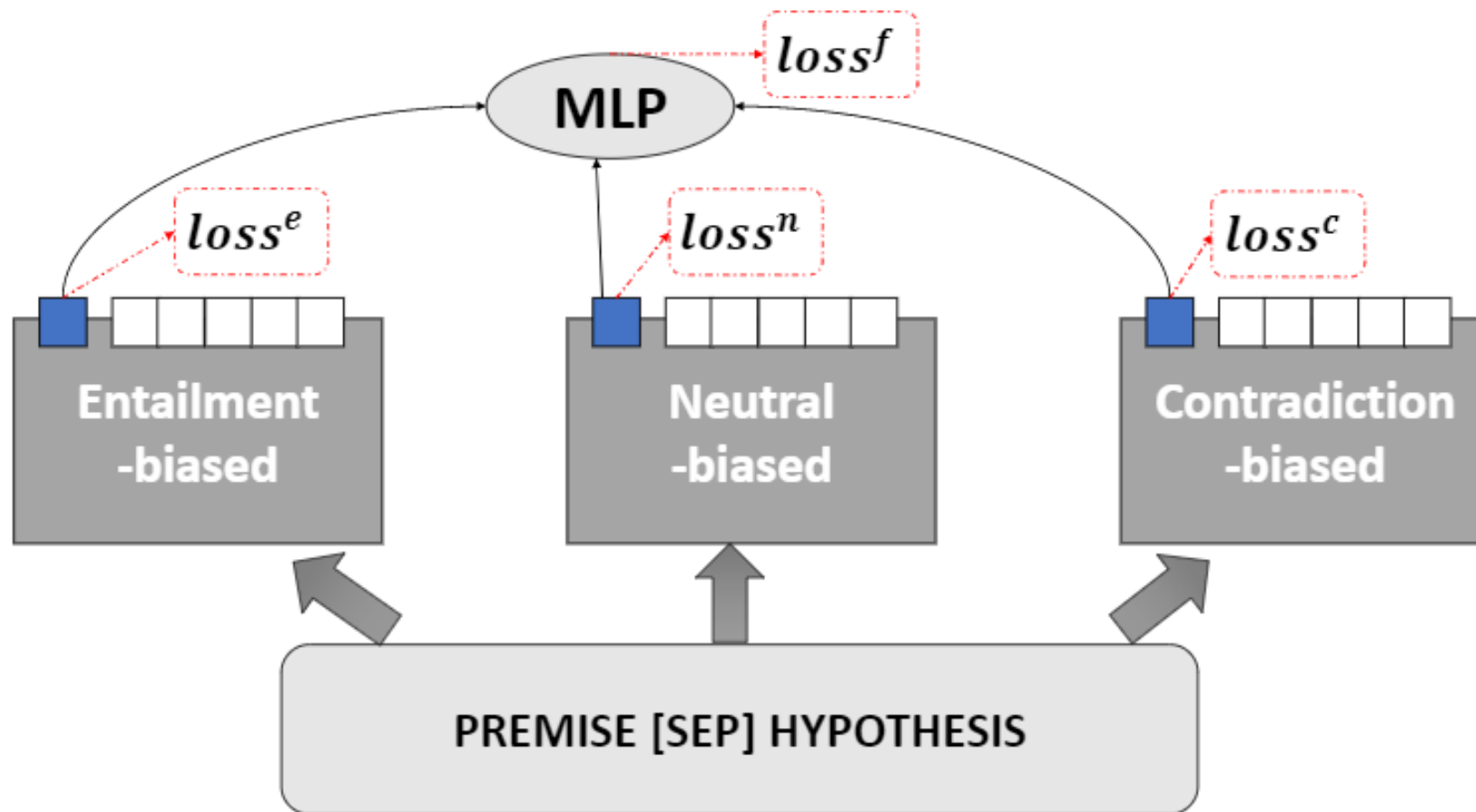
**Hypothesis:** that would be a good environment to work in



Disagreement

 Disagreement [ 2, 0, 0, 0, 0, -1, -2, -3 ]

# Model: Artificial Annotators (AAs)



$$P(y|\mathbf{x}) = \text{softmax}(\mathbf{W}_s \tanh(\mathbf{W}_t [\mathbf{e}; \mathbf{n}; \mathbf{c}]))$$

# AAs perform better across the board

	Dev		Test					
	Acc.	F1	Acc.	F1	Entail	Neutral	Contradict	Disagree
Always 0	55.00	39.03	45.42	28.37	0.00	0.00	0.00	62.46
CBOW	55.25	40.54	45.09	28.37	0.00	0.00	0.69	62.17
Heuristic	65.00	62.08	54.17	50.60	22.54	<b>52.94</b>	64.46	58.20
Vanilla BERT	63.71	63.54	62.50	61.93	59.26	49.64	69.09	61.93
Joint BERT	64.47	64.28	62.61	62.07	59.77	47.27	67.36	63.21
AAs (ours)	<b>65.15</b>	<b>64.41</b>	<b>65.60*</b>	<b>64.97*</b>	<b>61.07</b>	51.27	<b>70.89</b>	<b>66.49*</b>

Baselines and AAs overall performance on CB dev and test sets, and F1 scores of each class on the test set (average of 10 runs). \* indicates a statistically significant difference (t-test,  $p \leq 0.01$ ).



# AAs learn linguistic patterns and context-dependent inference better

Correct inference by Heuristic?	Correctly predicted (130)		Missed (110)	
	Acc.	F1	Acc.	F1
V. BERT	80.00	80.45	41.51	42.48
J. BERT	79.74	80.04	42.73	44.15
AAs	<b>84.37</b>	<b>84.85</b>	<b>46.97</b>	<b>48.75</b>

BERT-based models performance on test items correctly predicted by vs. items missed by linguistic rules.



# Error analysis

**Premise:** B: Yeah, and EDS is very particular about this, hair cuts, A: Wow. B: I mean it was like you can't have, you know, such and such facial hair, no beards, you know, and just really detailed. A: I don't know that **that would be a good environment to work in.**

**Hypothesis:** that would be a good environment to work in

Heuristics: C

V. BERT: C

J. BERT: D

AAs: C {C, C, C}



Disagreement [2, 0, 0, 0, 0, -1, -2, -3]

# Towards robust NLI

Our Artificial Annotators are a start in this direction but still far from succeeding (~ 66%).

A method which captures accurately the number of modes in the annotation distribution would lead to a better model.

# Thanks!



Code is available at:

<https://github.com/FrederickXZhang/FgNLI>



Contact: [zhang.9975@osu.edu](mailto:zhang.9975@osu.edu)