

# Efficient Argument Structure Extraction with Transfer Learning and Active Learning

Engineering

Bloomberg

## Findings of ACL 2022

Xinyu Hua, AI Researcher, Bloomberg  
Lu Wang, University of Michigan

**TechAtBloomberg.com**

© 2022 Bloomberg Finance L.P. All rights reserved.

# Roadmap

- Motivation
- Prior Work
- Task and Model
- Dataset
- Transfer Learning
- Active Learning
- Conclusion

# Motivation

## Comment:

I think this submission does not meet the community standard.

The originality of the approach is unclear. Most existing work (...) The difference here is (...) not meaningful.

Secondly, none of the baselines uses (...), which is unfair comparison.

Add

Comment



# Motivation

Support  
relations

## Comment:

I think this submission does not meet the community standard.

The originality of the approach is unclear.

Most existing work (...) The difference here is (...) not meaningful.

Secondly, none of the baselines uses (...), which is unfair comparison.

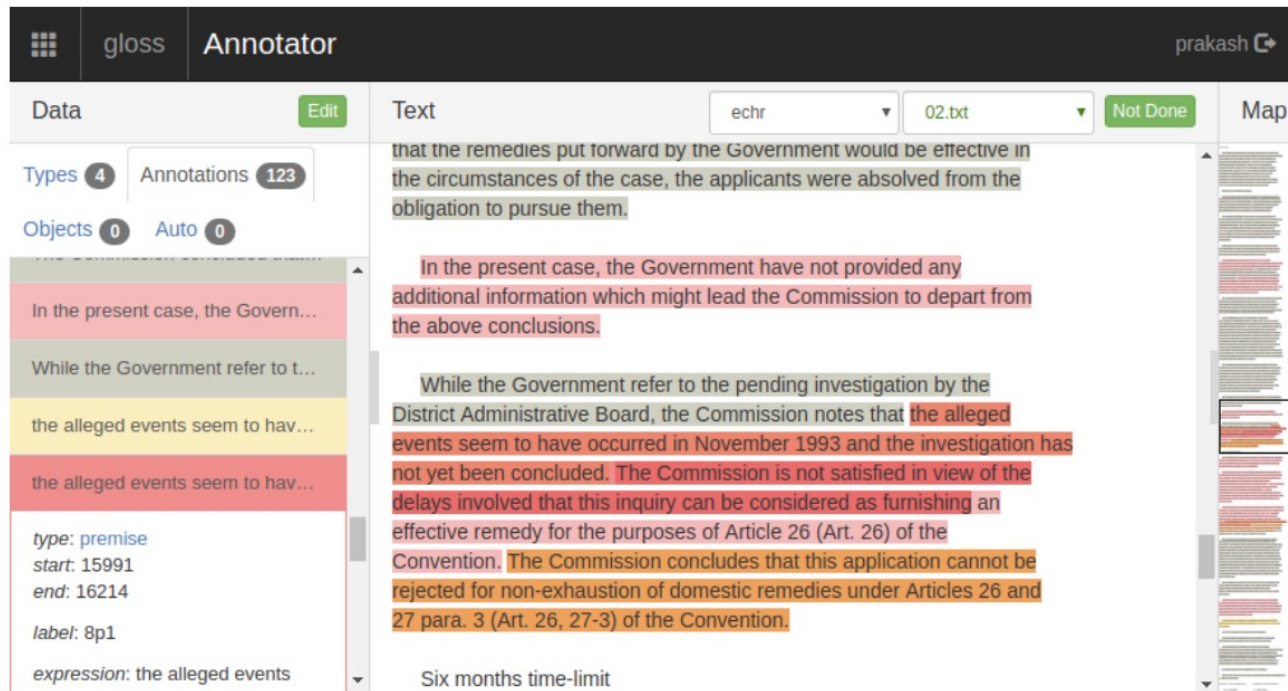
Add

Comment



# Motivation

- Annotating argument relations is difficult



Annotators (domain experts) need to scan through long documents.

Credit: Gloss interface, ECHR's Case Law [[Poudyal+, 2020](#)]

**TechAtBloomberg.com**

© 2022 Bloomberg Finance L.P. All rights reserved.

**Bloomberg**

Engineering

# Motivation

- Many existing (small) datasets exist, but no unified framework

*First, [cloning will be beneficial for many people who are in need of organ transplants]<sub>Claim2</sub>.  
[Cloned organs will match perfectly to the blood group and tissue of patients]<sub>Premise1</sub>  
since [they can be raised from cloned stem cells of the patient]<sub>Premise2</sub>. In addition, [it  
shortens the healing process]<sub>Premise3</sub>. Usually, [it is very rare to find an appropriate organ  
donor]<sub>Premise4</sub> and [by using cloning in order to raise required organs the waiting time  
can be shortened tremendously]<sub>Premise5</sub>.*

Student essays [Stab & Gurevych, 2017]

(1) \$400 is enough compensation,<sub>A</sub> as it can cover a one-way fare across the US.<sub>B</sub> I checked in a passenger on a \$98.00 fare from east coast to Las Vegas the other day.<sub>C</sub>

Online comments [Park & Cardie, 2018]

**Example 2** *[True acupuncture was associated with 0.8 fewer hot flashes per day than sham at 6 weeks,]<sub>1</sub> [but the difference did not reach statistical significance (95% CI, -0.7 to 2.4;  $p = 0.001$ )]*

Biomedical domain, paper abstract [Mayer+, 2020]

**“The notion of security of person has not been given an independent interpretation (see in this respect Selçuk and Asker v. Turkey, nos. 23184/94 and 23185/94, Commission’s report of 28 November 1996, §§ 185-187).”**

Legal domain, case law [Poudyal+, 2020]

TechAtBloomberg.com

© 2022 Bloomberg Finance L.P. All rights reserved.

Bloomberg

Engineering

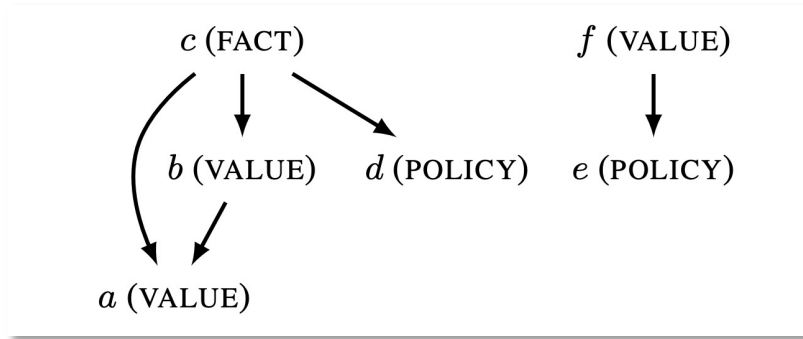
# Roadmap

- Motivation
- **Prior Work**
- Task and Model
- Dataset
- Transfer Learning
- Active Learning
- Conclusion

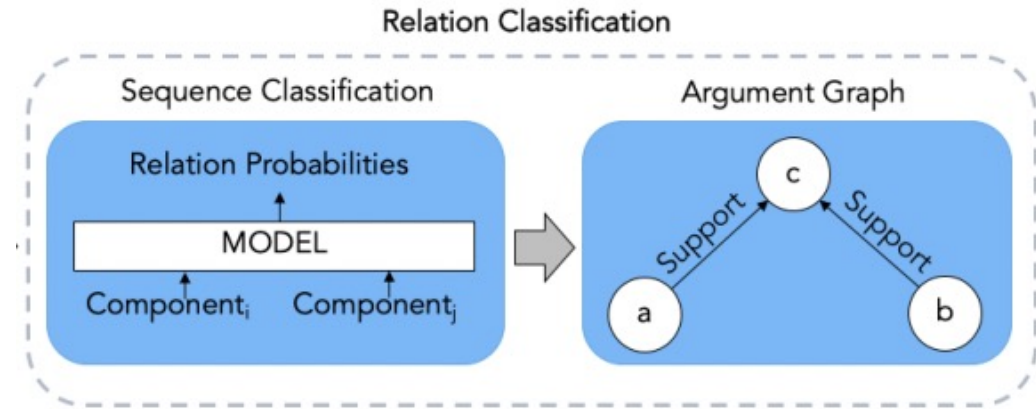


# Prior Work

- Argument Structure Prediction



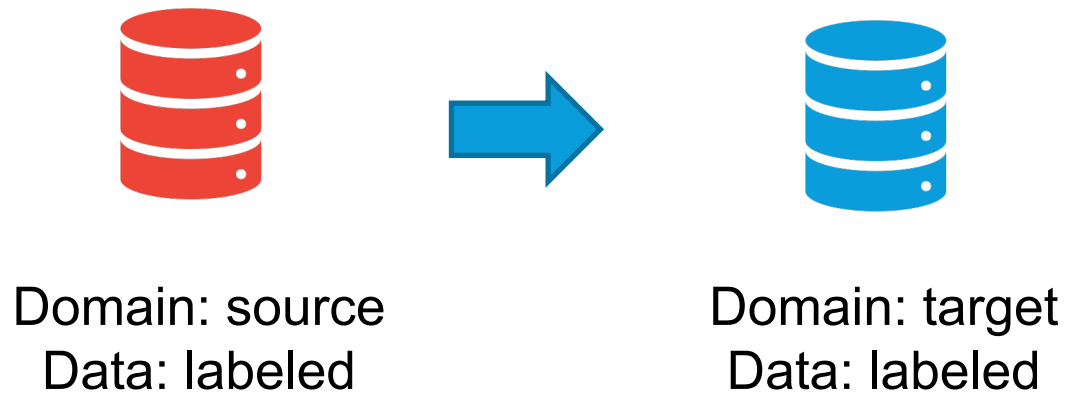
Factor graph with constraints  
[[Niculae, Park, and Cardie, 2017](#)]



Pairwise predictions  
[[Stab and Gurevych, 2017](#)]  
[[Mayer, Cabrio, and Villata, 2020](#)]

# Prior Work

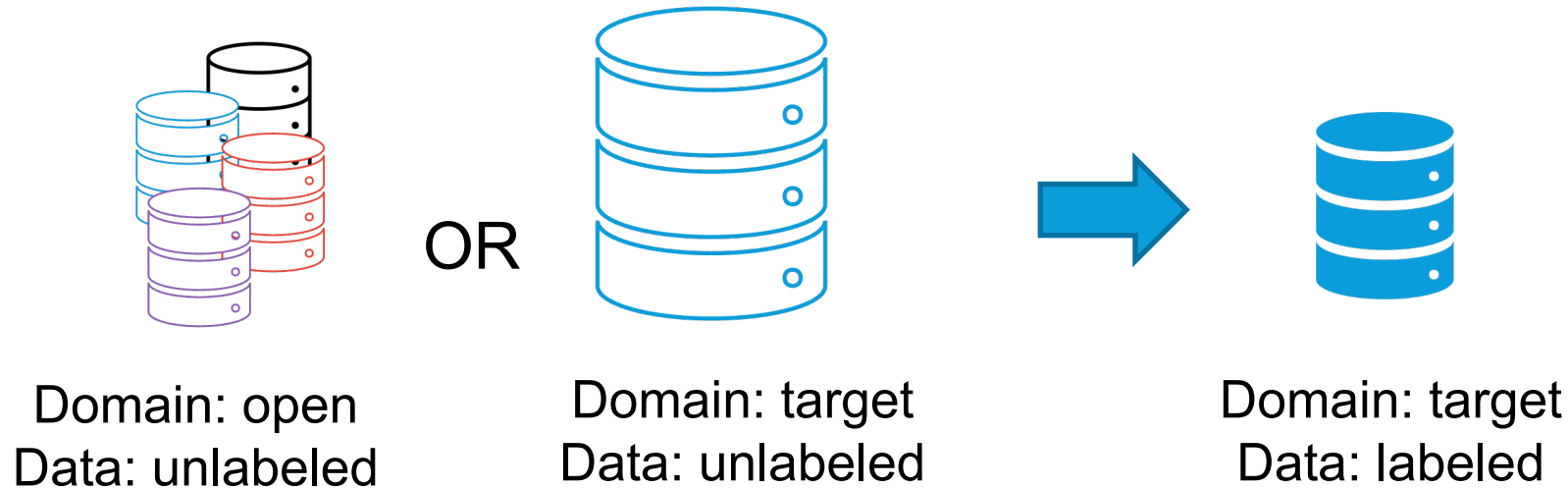
- Transfer Learning for Structured Prediction
  - Transductive
  - Inductive



Both source and target data are labeled (supervised training), on the **same task**

# Prior Work

- Transfer Learning for Structured Prediction
  - Transductive
  - Inductive

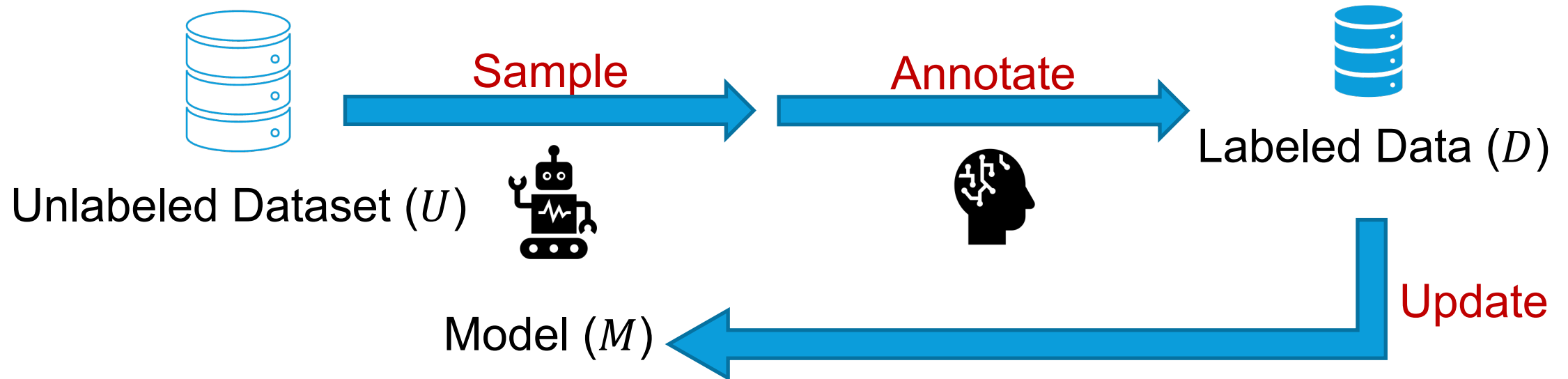


Transfer from unlabeled data, e.g., using self-supervised training



# Prior Work

- Active Learning [[Settles, 2009](#); [Aggarwal+, 2014](#)]
  - Unlabeled dataset is available
  - Annotation is subject to a budget
  - Goal is to select the most informative samples



# Roadmap

- Motivation
- Prior Work
- Task and Model
- Dataset
- Transfer Learning
- Active Learning
- Conclusion

# Task and Model

- Task Formulation

## Comment:

I think this submission does not meet the community standard.

The originality of the approach is unclear.

Most existing work (...) The difference here is (...) not meaningful.

Secondly, none of the baselines uses (...), which is unfair comparison.

Add [Comment](#)

## Propositions

$S_1$

$S_2$

$S_3$

$S_4$

$S_5$



# Task and Model

- Task Formulation

		To (head)				
		$s_1$	$s_2$	$s_3$	$s_4$	$s_5$
From (tail)	$s_1$					
	$s_2$	✓				
	$s_3$		✓			
	$s_4$					✓
	$s_5$	✓				

## Comment:

I think this submission does not meet the community standard.

The originality of the approach is unclear.

Most existing work (...) The difference here is (...) not meaningful.

Secondly, none of the baselines uses (...), which is unfair comparison.

Add

Comment

## Propositions

$s_1$

$s_2$

$s_3$

$s_4$

$s_5$

# Task and Model

- Task Formulation

From (tail) To (head)

	$s_1$	$s_2$	$s_3$	$s_4$	$s_5$
$s_1$					
$s_2$	✓				
$s_3$		✓			
$s_4$					✓
$s_5$	✓				

## Comment:

I think this submission does not meet the community standard.

The originality of the approach is unclear.

Most existing work (...) The difference here is (...) not meaningful.

Secondly, none of the baselines uses (...), which is unfair comparison.

Add

Comment

## Propositions

$s_1$

$s_2$

$s_3$

$s_4$

$s_5$

# Task and Model

- Task Formulation

From (tail)

	To (head)				
	$s_1$	$s_2$	$s_3$	$s_4$	$s_5$
$s_1$					
$s_2$	✓				
$s_3$		✓			
$s_4$					✓
$s_5$	✓				

**Comment:**

I think this submission does not meet the community standard.

The originality of the approach is unclear.  
Most existing work (...) The difference here is (...) not meaningful.

Secondly, none of the baselines uses (...), which is unfair comparison.

Add [Comment](#)

Propositions

$s_1$

$s_2$

$s_3$

$s_4$

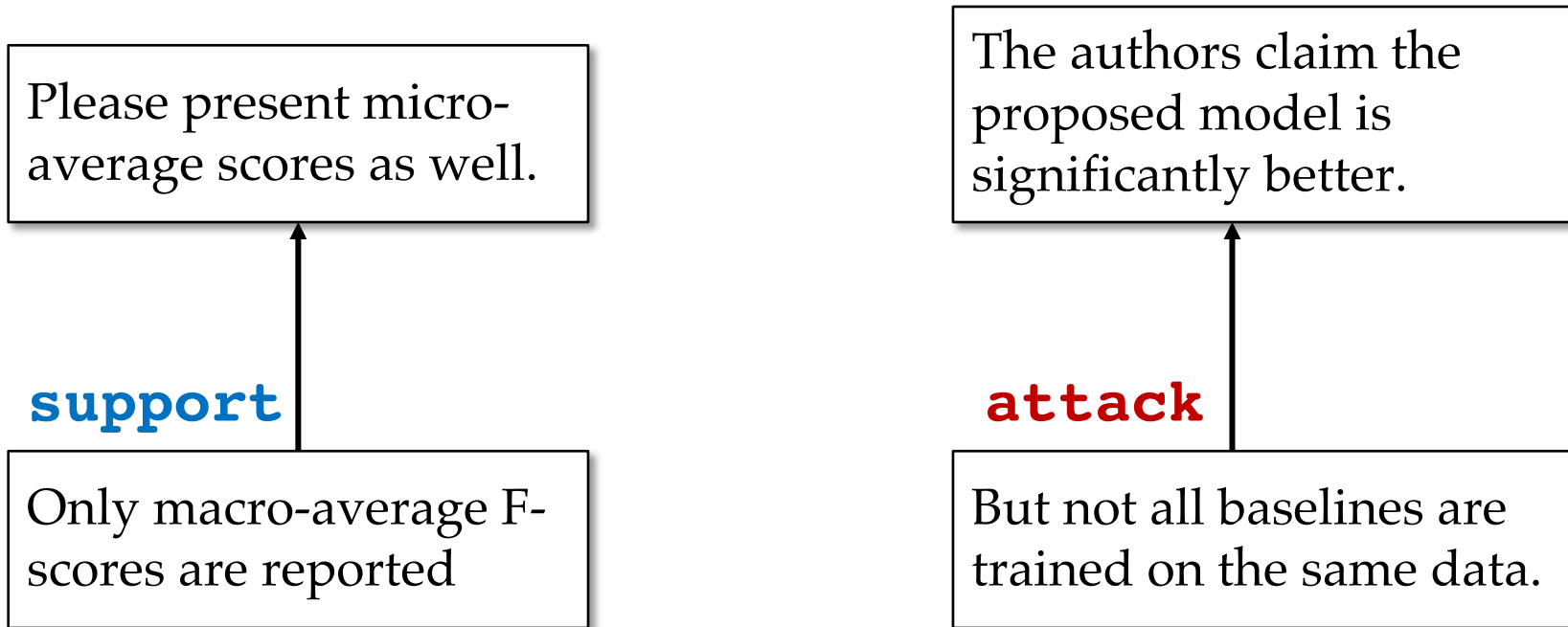
$s_5$

Simplified setting assumes head propositions ( $s_1, s_2, s_5$ ) are given.



# Task and Model

- Task Formulation
  - Types of relations: **support** and **attack**



# Task and Model

- Context-aware model

$S_j$  I think this submission does not meet the community standard.

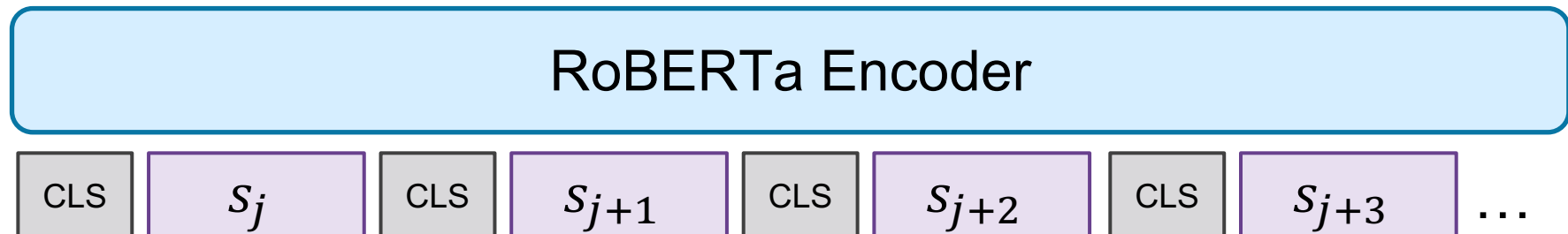
$S_{j+1}$  The originality of the approach is unclear.

$S_{j+2}$  Most existing work (...)

$S_{j+3}$  The difference here is (...) not meaningful.

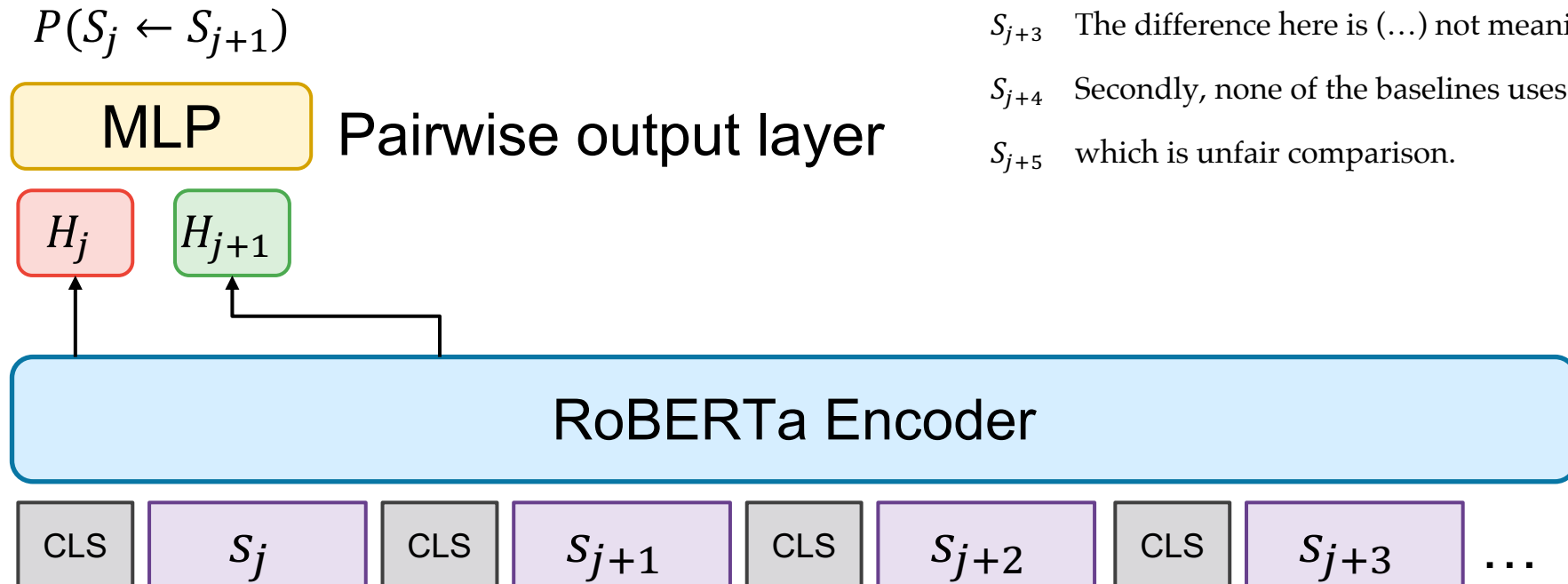
$S_{j+4}$  Secondly, none of the baselines uses (...),

$S_{j+5}$  which is unfair comparison.



# Task and Model

- Context-aware model



$S_j$  I think this submission does not meet the community standard.

$S_{j+1}$  The originality of the approach is unclear.

$S_{j+2}$  Most existing work (...)

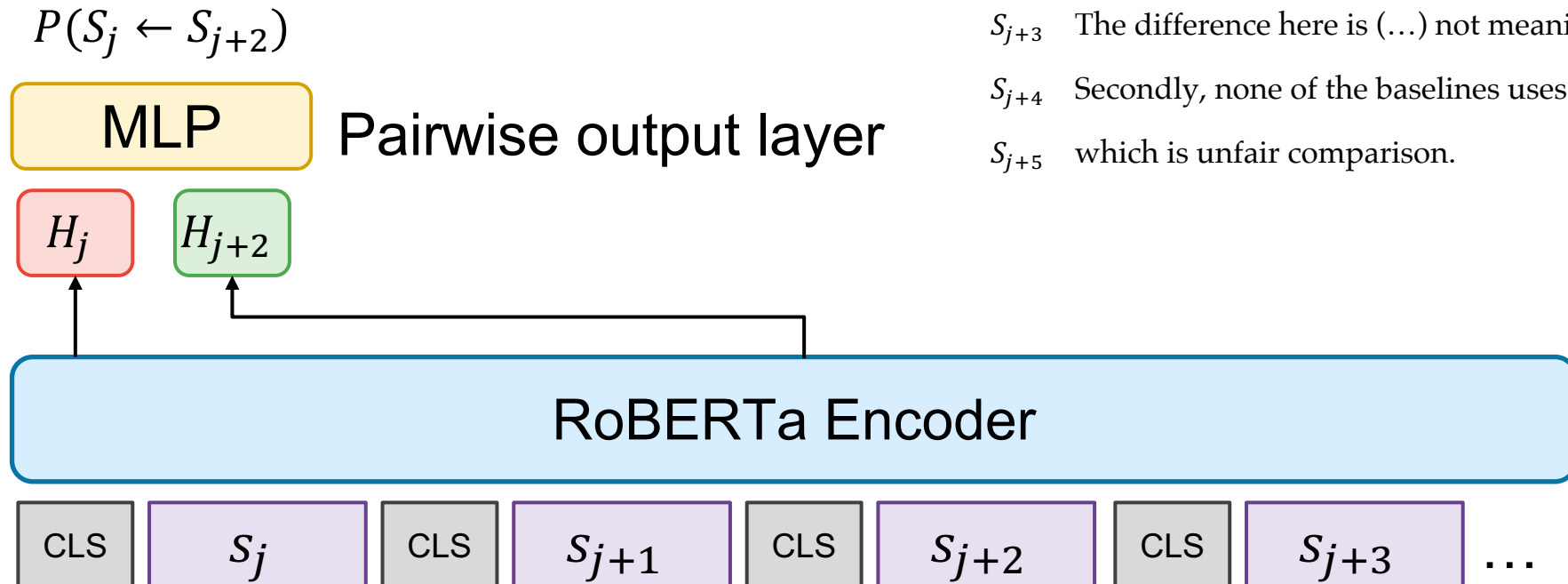
$S_{j+3}$  The difference here is (...) not meaningful.

$S_{j+4}$  Secondly, none of the baselines uses (...),

$S_{j+5}$  which is unfair comparison.

# Task and Model

- Context-aware model



$S_j$  I think this submission does not meet the community standard.

$S_{j+1}$  The originality of the approach is unclear.

$S_{j+2}$  Most existing work (...)

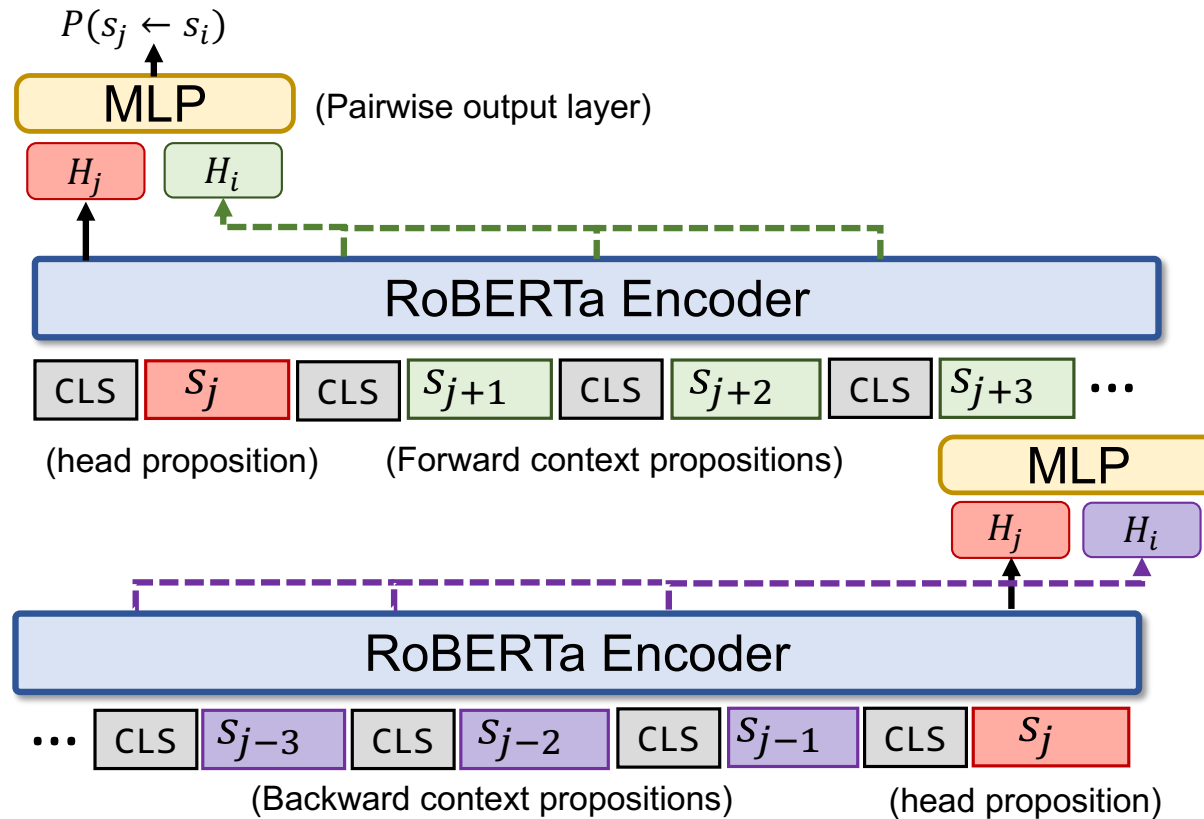
$S_{j+3}$  The difference here is (...) not meaningful.

$S_{j+4}$  Secondly, none of the baselines uses (...),

$S_{j+5}$  which is unfair comparison.

# Task and Model

- Context-aware model





# Roadmap

- Motivation
- Prior Work
- Task and Model
- **Dataset**
- Transfer Learning
- Active Learning
- Conclusion

# Dataset

- AMPERE++ (new annotation)
  - Domain: paper reviews from [openreview.net](https://openreview.net)
  - Originally collected in our prior work [Hua+, 2019]
  - 3,636 argument relations ([support](#) and [attack](#))
  - IAA: 0.654 (Fleiss' kappa)

# Dataset

- Essays [[Stab & Gurevych, 2017](#)]
- AbstRCT [[Mayer+, 2020](#)]
- ECHR [[Poudyal+, 2020](#)]
- CDCP [[Park & Cardie, 2018](#)]

# Dataset

- Essays [[Stab & Gurevych, 2017](#)]
- AbstRCT [Mayer+, 2020]
- ECHR [Poudyal+, 2020]
- CDCP [Park & Cardie, 2018]

*First, [cloning will be beneficial for many people who are in need of organ transplants]<sub>Claim2</sub>. [Cloned organs will match perfectly to the blood group and tissue of patients]<sub>Premise1</sub> since [they can be raised from cloned stem cells of the patient]<sub>Premise2</sub>. In addition, [it shortens the healing process]<sub>Premise3</sub>. Usually, [it is very rare to find an appropriate organ donor]<sub>Premise4</sub> and [by using cloning in order to raise required organs the waiting time can be shortened tremendously]<sub>Premise5</sub>.*

# Dataset

- Essays [Stab & Gurevych, 2017]
- AbstRCT [[Mayer+, 2020](#)]
- ECHR [Poudyal+, 2020]
- CDCP [Park & Cardie, 2018]

**Example 2** *[True acupuncture was associated with 0.8 fewer hot flashes per day than sham at 6 weeks,]<sub>1</sub> [but the difference did not reach statistical significance (95% CI, -0.7 to 2.4; P = .3).]<sub>2</sub>*



# Dataset

- Essays [Stab & Gurevych, 2017]
- AbstRCT [Mayer+, 2020]
- ECHR [[Poudyal+, 2020](#)]
- CDCP [Park & Cardie, 2018]

**“The notion of security of person has not been given an independent interpretation (*see in this respect Selçuk and Asker v. Turkey, nos. 23184/94 and 23185/94, Commission’s report of 28 November 1996, §§ 185-187*).”**

# Dataset

- Essays [Stab & Gurevych, 2017]
- AbstRCT [Mayer+, 2020]
- ECHR [Poudyal+, 2020]
- CDCP [[Park & Cardie, 2018](#)]

(1) \$400 is enough compensation,<sub>A</sub> as it can cover a one-way fare across the US.<sub>B</sub> I checked in a passenger on a \$98.00 fare from east coast to Las Vegas the other day.<sub>C</sub>

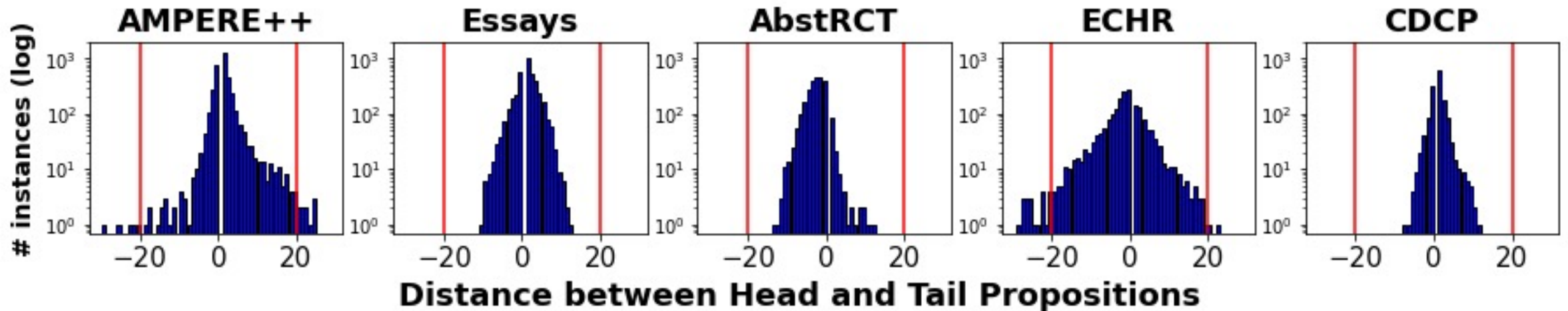
# Dataset

- Statistics

	AMPERE++	Essays	AbstRCT	ECHR	CDCP
# Documents	400	402	700	42	731
# Propositions	10.4K	12.4K	5.7K	6.3K	4.9K
# Support Rel.	3,370	3,613	2,402	1,946	1,426
# Attack Rel.	266	219	70	0	0
# Head Prop.	2,268	1,707	1,138	741	1,037
Density	22%	14%	20%	12%	21%

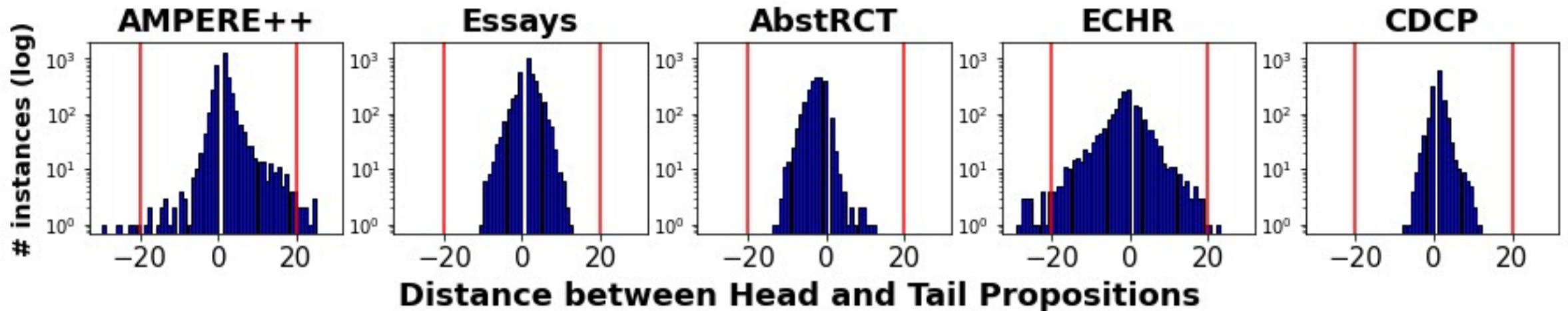
# Dataset

- Statistics: distribution of head-tail distance



# Dataset

- Statistics: distribution of head-tail distance



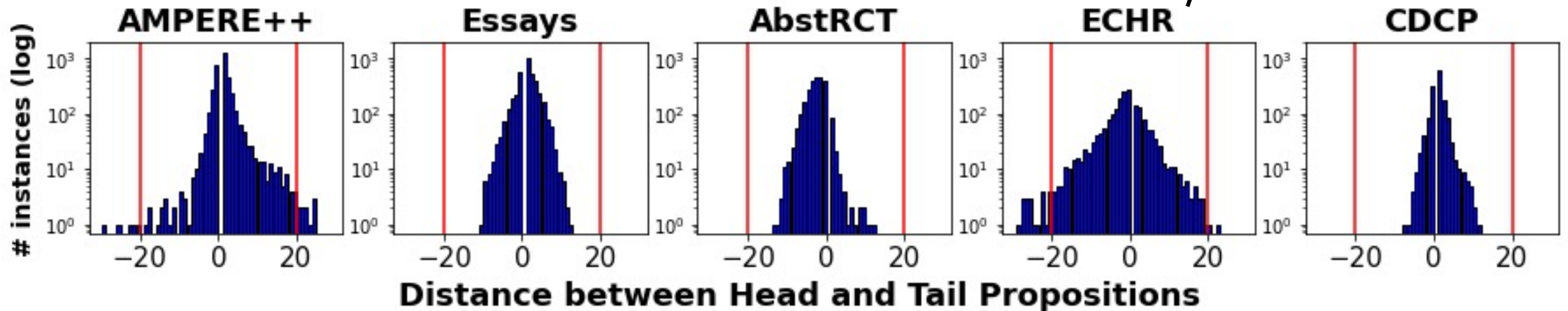
Most relations span less than 20 propositions.



# Dataset

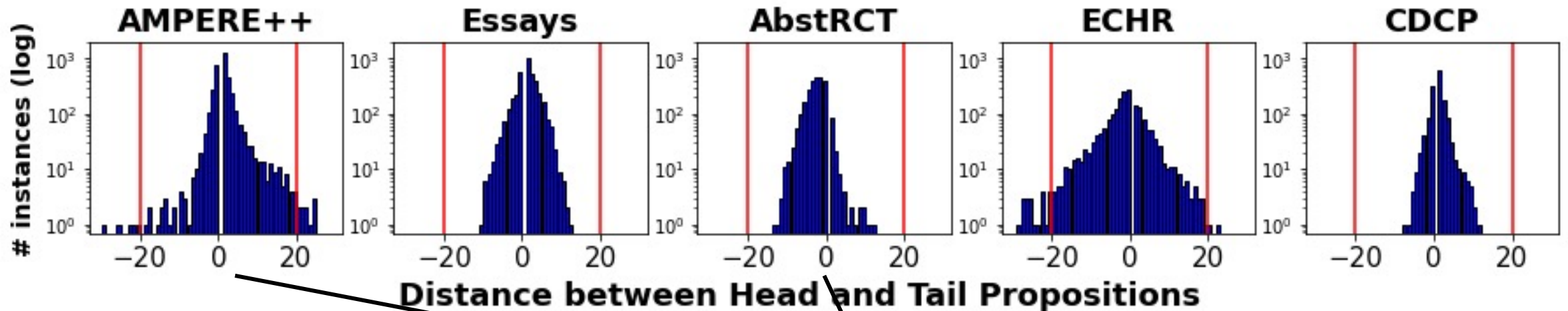
- Statistics: distribution of head-tail distance

Long-term relations are more common in the legal domain



# Dataset

- Statistics: distribution of head-tail distance



Review: verdicts first, support later  
Paper abstract: evidence first, conclusion later

**TechAtBloomberg.com**

© 2022 Bloomberg Finance L.P. All rights reserved.

**Bloomberg**

Engineering

# Standard Supervised Setting

- Baselines and our context-aware model
- Macro F1 scores

	AMPERE++	Essays	AbstRCT	ECHR	CDCP
SVM-linear	24.82	28.69	33.60	21.18	29.01
SVM-RBF	26.38	31.68	32.65	21.36	30.34
SEQPAIR	23.40	38.37	<b>66.96</b>	13.76	35.23
OURS (head given)	<b>77.64</b>	<b>71.30</b>	63.62	<b>70.82</b>	<b>70.37</b>
OURS (end-to-end)	74.34	67.68	63.73	61.35	63.13

# Standard Supervised Setting

- Baselines and our context-aware model
- Macro F1 scores

	AMPERE++	Essays	AbstRCT	ECHR	CDCP
SVM-linear	24.82	28.69	33.60	21.18	29.01
SVM-RBF	26.38	31.68	32.65	21.36	30.34
SEQPAIR	23.40	38.37	<b>66.96</b>	13.76	35.23
OURS (head given)	<b>77.64</b>	<b>71.30</b>	63.62	<b>70.82</b>	<b>70.37</b>
OURS (end-to-end)	74.34	67.68	63.73	61.35	63.13

Takeaways:

- 1) Context-aware model is generally much better
- 2) End-to-end and simplified (head given) setting are close

# Roadmap

- Motivation
- Prior Work
- Task and Model
- Dataset
- **Transfer Learning**
- Active Learning
- Conclusion

# Transfer Learning

- Transductive TL
  - Same task, different domains (datasets)
  - (Source) model weights as (target) initialization



# Transfer Learning

- Transductive TL
  - Same task, different domains (datasets)
  - (Source) model weights as (target) initialization

		Target Domain				
Source Domain		AMPERE++	Essays	AbstRCT	ECHR	CDCP
	AMPERE++		73.84	63.42	76.50	75.93
	Essays	77.93		60.62	68.72	74.11
	AbstRCT	76.29	71.17		73.31	69.17
	ECHR	77.69	70.82	47.91		69.30
	CDCP	77.87	68.37	62.38	72.03	

# Transfer Learning

- Transductive TL
  - Same task, different domains (datasets)
  - (Source) model weights as (target) initialization

Transfer settings that outperform standard supervised setting are **highlighted**

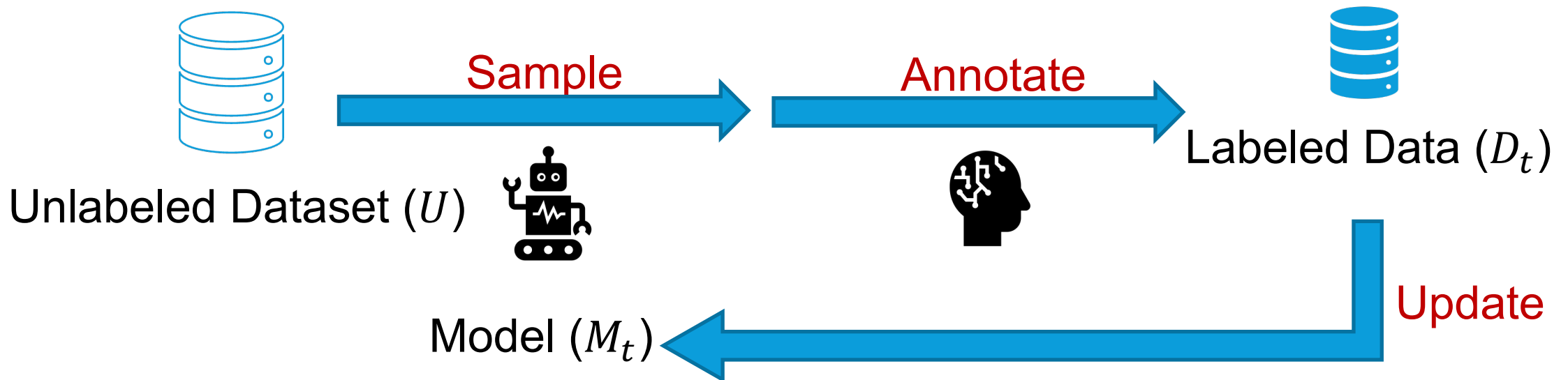
		Target Domain				
		AMPERE++	Essays	AbstRCT	ECHR	CDCP
Source Domain	AMPERE++		73.84	63.42	76.50	75.93
	Essays	77.93		60.62	68.72	74.11
	AbstRCT	76.29	71.17		73.31	69.17
	ECHR	77.69	70.82	47.91		69.30
	CDCP	77.87	68.37	62.38	72.03	

# Roadmap

- Motivation
- Prior Work
- Task and Model
- Dataset
- Transfer Learning
- **Active Learning**
- Conclusion

# Active Learning

- Experiment settings:
  - 10 iterations, 500 samples per iteration



# Active Learning

- Acquisition strategies
  - RANDOM
  - MAX-ENTROPY [[Lewis & Gale, 1994](#); [Joshi+, 2009](#)]
  - BALD [[Houlsby+, 2011](#)]
  - CORESET [[Sener & Savarese, 2018](#)]

# Active Learning

- Acquisition strategies
    - RANDOM
    - MAX-ENTROPY [[Lewis & Gale, 1994](#); [Joshi+, 2009](#)]
    - BALD [[Houlsby+, 2011](#)]
    - CORESET [[Sener & Savarese, 2018](#)]
- } Picks the most uncertain samples
- Maximizes sample diversity

# Active Learning

- Acquisition strategies
  - RANDOM
  - MAX-ENTROPY [Lewis & Gale, 1994; Joshi+, 2009]
  - BALD [Houlsby+, 2011]
  - CORESET [Sener & Savarese, 2018]
  - NOVEL-VOCAB

$$\text{novelty-score}(\mathbf{s}_i) = \sum_{w_t \in \mathbf{s}_i} \frac{f_{i,t}}{(1 + \mathcal{V}(w_t))}$$

Frequency of word  $w_t$  in sample  $S_i$

Frequency of word  $w_t$  in the labeled  $D_t$



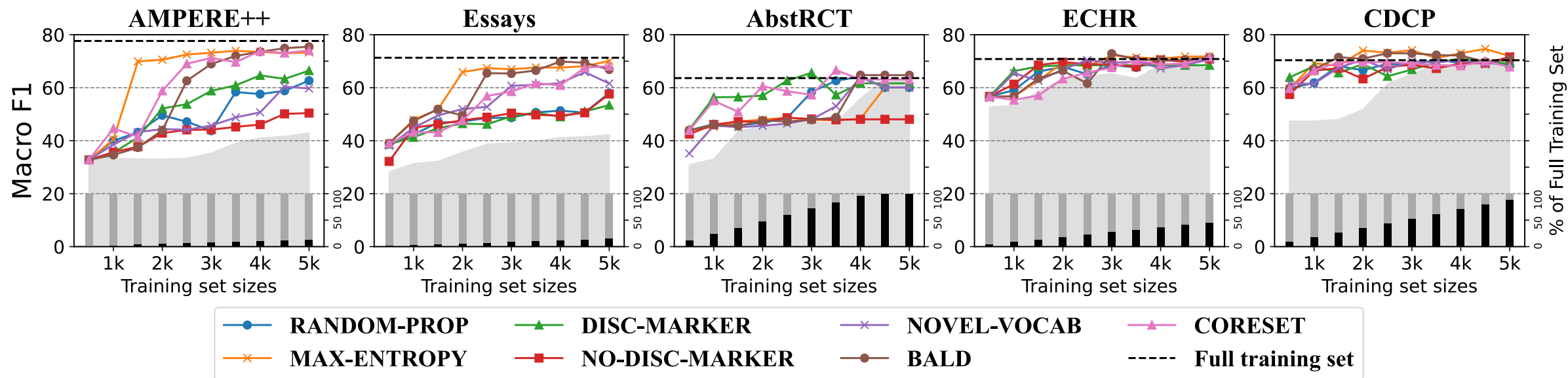
# Active Learning

- Acquisition strategies
  - RANDOM
  - MAX-ENTROPY [[Lewis & Gale, 1994](#); [Joshi+, 2009](#)]
  - BALD [[Houlsby+, 2011](#)]
  - CORESET [[Sener & Savarese, 2018](#)]
  - NOVEL-VOCAB
  - DISC-MARKER

because	therefore	however
although	though	nevertheless
nonetheless	thus	hence
consequently	for this reason	due to
in particular	particularly	specifically
in fact	actually	but

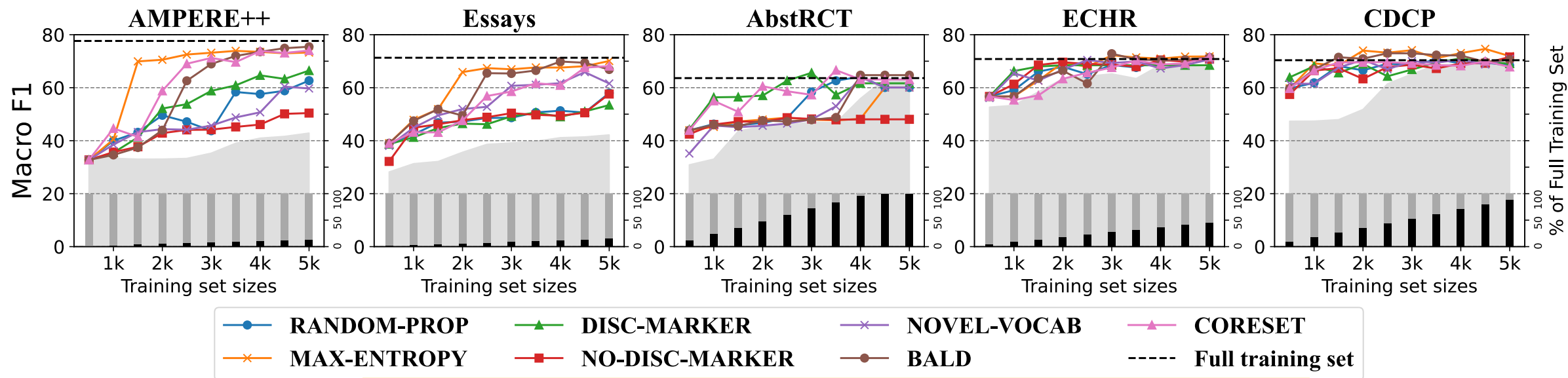
# Active Learning

- Results



# Active Learning

- Results



## Takeaways:

- 1) AL helps more in the early stage (low-resource).
- 2) Model-independent strategies achieve competitive performance, yet they are much faster.

# Conclusion

- We present a simple yet effective framework for argument structure extraction.
- We release AMPERE++, a newly annotated dataset on peer reviews.
- We showcase two data efficient learning methods (transfer learning and active learning) using our model.



# Questions?

**Bloomberg**  
**Engineering**



<https://xinyuhua.github.io/Resources/acl22/>



<https://zenodo.org/record/6362430>



[xhua22@bloomberg.net](mailto:xhua22@bloomberg.net)

**TechAtBloomberg.com**

© 2022 Bloomberg Finance L.P. All rights reserved.

## Transfer Learning

- Inductive TL
  - Same domain, different tasks (self-supervision)
  - MLM: masked language model
  - Context-Pert: context-aware sentence perturbation

	AMPERE++	Essays	AbstRCT
MLM	78.10	74.21	64.48
Context-Pert	79.01	68.36	59.47

## Active Learning

- Acquisition strategies
  - RANDOM: RANDOM-CTX vs. RANDOM-PROP
  - MAX-ENTROPY [Lewis & Gale, 1994; Joshi+, 2009]
  - BALD [Houlsby+, 2011]
  - CORESET [Sener & Savarese, 2018]

RANDOM-CTX



RANDOM-PROP



S<sub>1</sub>

S<sub>2</sub>

S<sub>3</sub>

S<sub>4</sub>

S<sub>5</sub>

S<sub>6</sub>

S<sub>7</sub>

S<sub>8</sub>

S<sub>9</sub>

S<sub>10</sub>

S<sub>11</sub>