

Adapting Coreference Resolution Models through Active Learning

ACL 2022

**Michelle Yuan¹ Patrick Xia² Chandler May²
Benjamin Van Durme² Jordan Boyd-Graber¹**

¹Department of Computer Science
University of Maryland

²Department of Computer Science
Johns Hopkins University



Coreference Resolution (CR)

The task of discovering spans of text that refer to the same entity

Source
(Finance)

Traders said **municipals** were underpinned by influences, including the climb in Treasury issue prices. Also, **municipal bonds** lured buying because the stock market remains wobbly, traders contended. Mainly though, it was a favorable outlook for yesterday's new supply that propped up **municipals**, some traders said. Among the new issues was Massachusetts's \$230 million of **general obligation bonds**. **The bonds** were won by a Goldman Sachs & Co. group with a true interest cost of 7.17%.

Cluster: { **municipals**, **municipal bonds**, **municipals**, **general obligation bonds**, **The bonds** }

Coreference Resolution (CR)

The task of discovering spans of text that refer to the same entity

Source
(Finance)

Traders said **municipals** were underpinned by influences, including the climb in Treasury issue prices. Also, **municipal bonds** lured buying because the stock market remains wobbly, traders contended. Mainly though, it was a favorable outlook for yesterday's new supply that propped up **municipals**, some traders said. Among the new issues was Massachusetts's \$230 million of **general obligation bonds**. **The bonds** were won by a Goldman Sachs & Co. group with a true interest cost of 7.17%.

Cluster: { **municipals** , **municipal bonds** , **municipals** , **general obligation bonds** , **The bonds** }

Neural, end-to-end models (Lee et al., 2018; Joshi et al., 2020) are SOTA for OntoNotes 5.0

Problem: Adapting CR Models

Models trained on OntoNotes may not immediately generalize to new domains

Source
(Finance)

Traders said **municipals** were underpinned by influences, including the climb in Treasury issue prices. Also, **municipal bonds** lured buying because the stock market remains wobbly, traders contended. Mainly though, it was a favorable outlook for yesterday's new supply that propped up **municipals**, some traders said. Among the new issues was Massachusetts's \$230 million of **general obligation bonds**. **The bonds** were won by a Goldman Sachs & Co. group with a true interest cost of 7.17%.

Cluster: { **municipals** , **municipal bonds** , **municipals** , **general obligation bonds** , **The bonds** }

Target
(Science)

A molecule is a group of atoms held together by **chemical bonds**. Imagine you and your friends standing in a circle holding hands. Each person stands for one atom, your hands represent **the bonds**, and the entire circle represents a molecule.

Problem: Adapting CR Models

Models trained on OntoNotes may not immediately generalize to new domains

Source
(Finance)

Traders said **municipals** were underpinned by influences, including the climb in Treasury issue prices. Also, **municipal bonds** lured buying because the stock market remains wobbly, traders contended. Mainly though, it was a favorable outlook for yesterday's new supply that propped up **municipals**, some traders said. Among the new issues was Massachusetts's \$230 million of **general obligation bonds**. **The bonds** were won by a Goldman Sachs & Co. group with a true interest cost of 7.17%.

Cluster: { **municipals**, **municipal bonds**, **municipals**, **general obligation bonds**, **The bonds** }

Target
(Science)

A molecule is a group of atoms held together by **chemical bonds**. Imagine you and your friends standing in a circle holding hands. Each person stands for one atom, your hands represent **the bonds**, and the entire circle represents a molecule.

Impedes immediate application for scenarios like distinguishing entities in scientific articles

Problem: Adapting CR Models

- (Xia and Van Durme, 2021) show the benefits of *continued training* where a model trained on OntoNotes is further trained on the target dataset

Problem: Adapting CR Models

- (Xia and Van Durme, 2021) show the benefits of *continued training* where a model trained on OntoNotes is further trained on the target dataset
- However, they assume labeled data already exist in the target domain

Problem: Adapting CR Models

- (Xia and Van Durme, 2021) show the benefits of *continued training* where a model trained on OntoNotes is further trained on the target dataset
- However, they assume labeled data already exist in the target domain
- How can we adapt CR models without requiring large amounts of newly annotated data?

Method: Active Learning

- Use active learning to find particular spans of text for users to label

Method: Active Learning

- Use active learning to find particular spans of text for users to label
- The goal is to adapt the model to the target domain by continue training it on spans labeled from active learning

What should we label?

A fantastic experience, very informative, very time consuming but enjoyable. So much information to take in about Guinness that you would've never known. For example, the brewery hired the statistician Willam Gosset in 1899. The "student" was known for developing the Student's t-test, a well-known method in statistical inference.

What should we label?

A fantastic experience, very informative, very time consuming but enjoyable. So much information to take in about Guinness that you would've never known. For example, the brewery hired the statistician Willam Gosset in 1899. The "student" was known for developing the Student's t-test, a well-known method in statistical inference.

Uncertainty in mention detection

What should we label?

A fantastic experience, very informative, very time consuming but enjoyable. So much information to take in about Guinness that you would've never known. For example, the brewery hired the statistician Willam Gosset in 1899. The "student" was known for developing the Student's t-test, a well-known method in statistical inference.

Uncertainty in mention clustering

What should we label?

A fantastic experience, very informative, very time consuming but enjoyable. So much information to take in about Guinness that you would've never known. For example, the brewery hired the statistician Willam Gosset in 1899. The "student" was known for developing the Student's t-test, a well-known method in statistical inference.

Uncertainty in mention clustering conditioned on
mention detection

What should we label?

A fantastic experience, very informative, very time consuming but enjoyable. So much information to take in about Guinness that you would've never known. For example, the brewery hired the statistician Willam Gosset in 1899. The “student” was known for developing the Student’s t-test, a well-known method in statistical inference.

Uncertainty in both mention detection and mention clustering

Experiments: Strategies

1. **ment-ent:** Mention detection entropy
2. **clust-ent:** Mention clustering entropy
3. **cond-ent:** Conditional entropy
4. **joint-ent** Joint entropy

Experiments: Strategies

1. **ment-ent:** Mention detection entropy
2. **clust-ent:** Mention clustering entropy
3. **cond-ent:** Conditional entropy
4. **joint-ent** Joint entropy
5. **random:** Randomly sample from all spans in the document

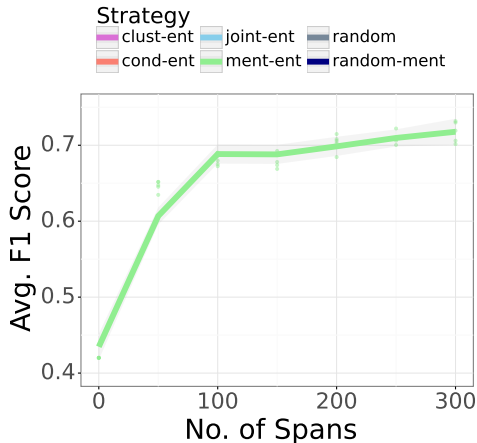
Experiments: Strategies

1. **ment-ent:** Mention detection entropy
2. **clust-ent:** Mention clustering entropy
3. **cond-ent:** Conditional entropy
4. **joint-ent** Joint entropy
5. **random:** Randomly sample from all spans in the document
6. **random-ment:** Randomly sample only from the pool of spans that are likely entity mentions

Experiments: Datasets

1. **OntoNotes 5.0 (source):** Most common dataset for training and evaluating CR that contains news articles and telephone conversations (Pradhan et al., 2013). Only non-singletons are annotated.
2. **PreCo (target):** Large corpus of grade-school reading comprehension texts with annotated singletons (Chen et al., 2018).

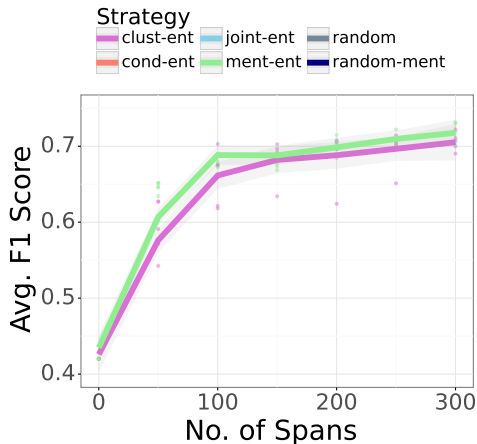
Comparing Strategies



- Test Avg. F1 on PreCo
- For each cycle, we simulate labeling fifty spans from one document

Mention Detection Entropy

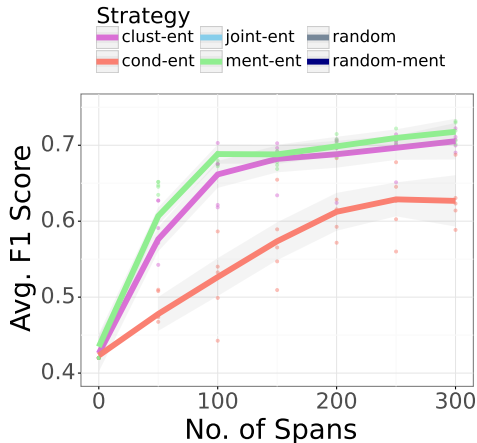
Comparing Strategies



- Test Avg. F1 on PreCo
- For each cycle, we simulate labeling fifty spans from one document

Mention Clustering Entropy

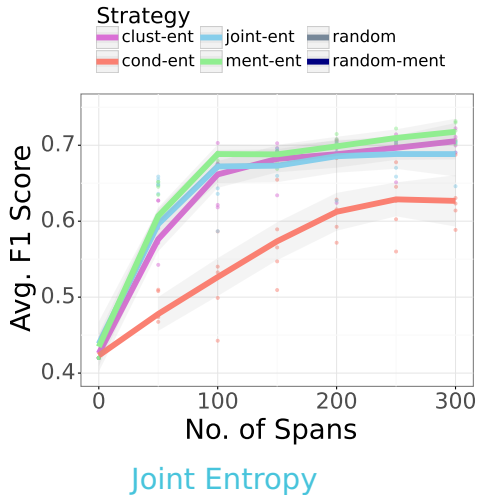
Comparing Strategies



Conditional Entropy

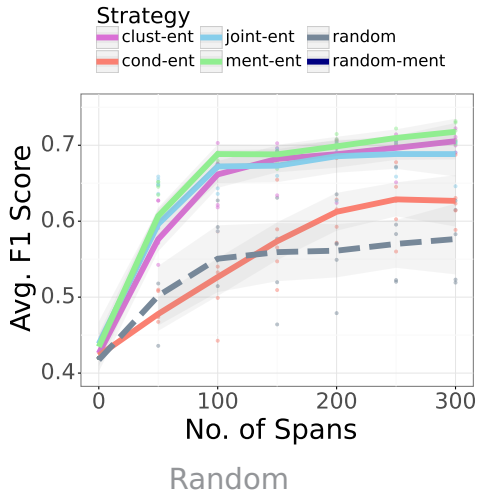
- Test Avg. F1 on PreCo
- For each cycle, we simulate labeling fifty spans from one document

Comparing Strategies



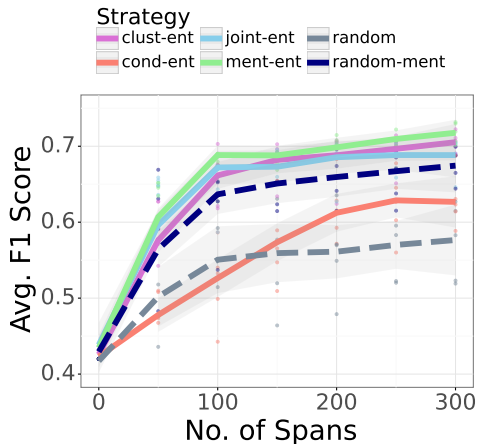
- Test Avg. F1 on PreCo
- For each cycle, we simulate labeling fifty spans from one document

Comparing Strategies



- Test Avg. F1 on PreCo
- For each cycle, we simulate labeling fifty spans from one document

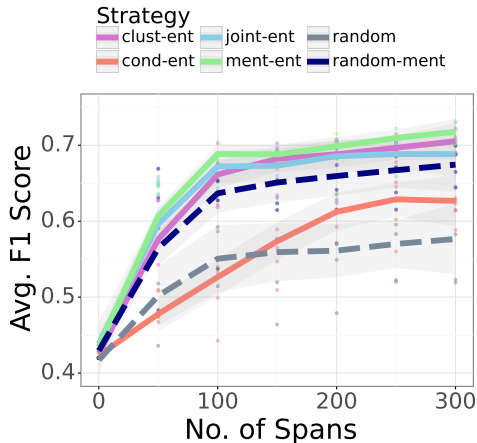
Comparing Strategies



Random Entity Mentions

- Test Avg. F1 on PreCo
- For each cycle, we simulate labeling fifty spans from one document

Comparing Strategies



- Test Avg. F1 on PreCo
- For each cycle, we simulate labeling fifty spans from one document
- **Ment-ent**, **clust-ent**, and **joint-ent** are effective while **random** performs worst

Should we label spans within or across documents?

A fantastic experience, very informative, very time consuming but enjoyable. So much information to take in about Guinness that you would've never known. For example, the brewery hired the statistician Willam Gosset in 1899. The "student" was known for developing the Student's t-test, a well-known method in statistical inference.

Ha'penny Bridge might look like it's just another bridge. But if you read about the history, you will know how significant this bridge is. Apparently half a penny was the toll that had to be paid from 1816 until the year 1919 in order to cross the Liffey Bridge.

Lovely park. Easy to get to on public transport. I recommend getting the bikes for hire when you get there, it made getting around really easy and you can cut across the fields to go and see the deer more easily!

Should we label spans within or across documents?

A fantastic experience, very informative, very time consuming but enjoyable. So much information to take in about Guinness that you would've never known. For example, the brewery hired the statistician Willam Gosset in 1899. The "student" was known for developing the Student's t-test, a well-known method in statistical inference.

Ha'penny Bridge might look like it's just another bridge. But if you read about the history, you will know how significant this bridge is. Apparently half a penny was the toll that had to be paid from 1816 until the year 1919 in order to cross the Liffey Bridge.

Lovely park. Easy to get to on public transport. I recommend getting the bikes for hire when you get there, it made getting around really easy and you can cut across the fields to go and see the deer more easily!

Within

Should we label spans within or across documents?

A fantastic experience, very informative, very time consuming but enjoyable. So much information to take in about Guinness that you would've never known. For example, the brewery hired the statistician Willam Gosset in 1899. The "student" was known for developing the Student's t-test, a well-known method in statistical inference.

Ha'penny Bridge might look like it's just another bridge. But if you read about the history, you will know how significant this bridge is. Apparently half a penny was the toll that had to be paid from 1816 until the year 1919 in order to cross the Liffey Bridge.

Lovely park. Easy to get to on public transport. I recommend getting the bikes for hire when you get there, it made getting around really easy and you can cut across the fields to go and see the deer more easily!

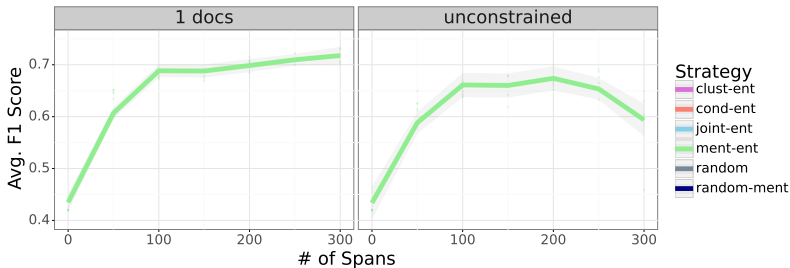
Across

Varying Number of Documents Read

- Test Avg. F1 on PreCo of each strategy
- On each cycle, sample fifty spans from either one document or across many documents

Varying Number of Documents Read

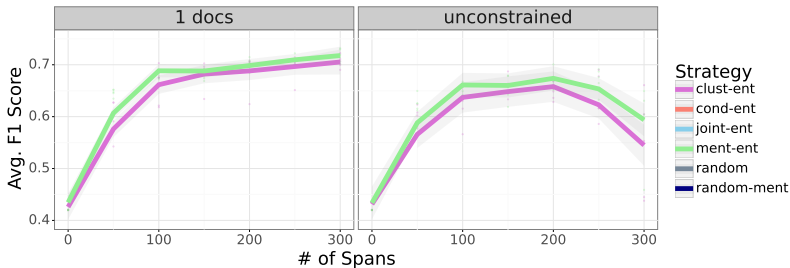
- Test Avg. F1 on PreCo of each strategy
- On each cycle, sample fifty spans from either one document or across many documents



Mention Detection Entropy

Varying Number of Documents Read

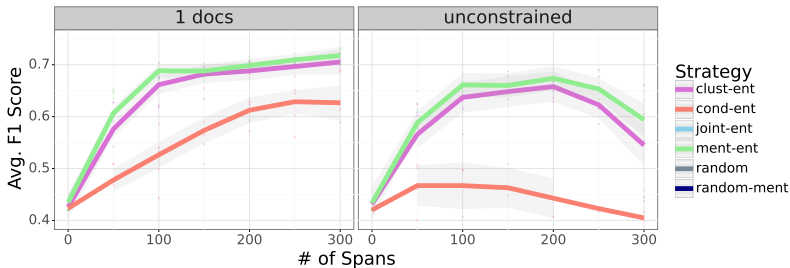
- Test Avg. F1 on PreCo of each strategy
- On each cycle, sample fifty spans from either one document or across many documents



Mention Clustering Entropy

Varying Number of Documents Read

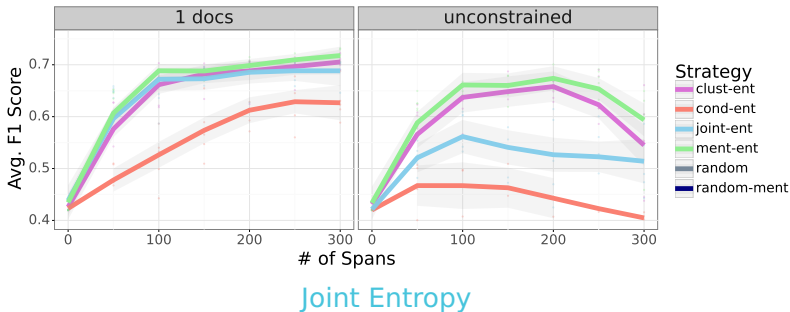
- Test Avg. F1 on PreCo of each strategy
- On each cycle, sample fifty spans from either one document or across many documents



Conditional Entropy

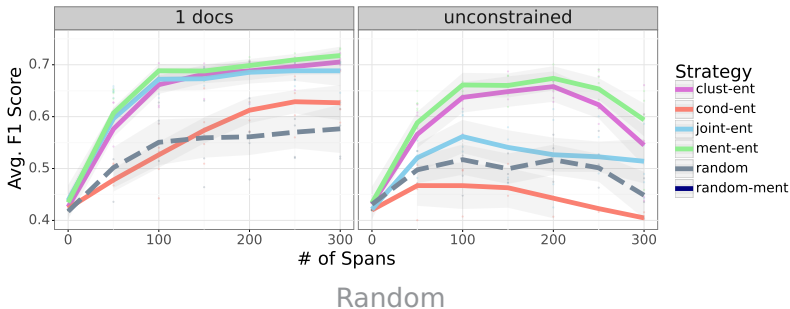
Varying Number of Documents Read

- Test Avg. F1 on PreCo of each strategy
- On each cycle, sample fifty spans from either one document or across many documents



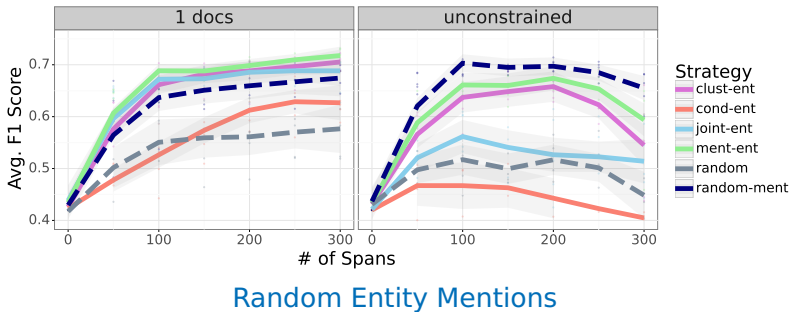
Varying Number of Documents Read

- Test Avg. F1 on PreCo of each strategy
- On each cycle, sample fifty spans from either one document or across many documents



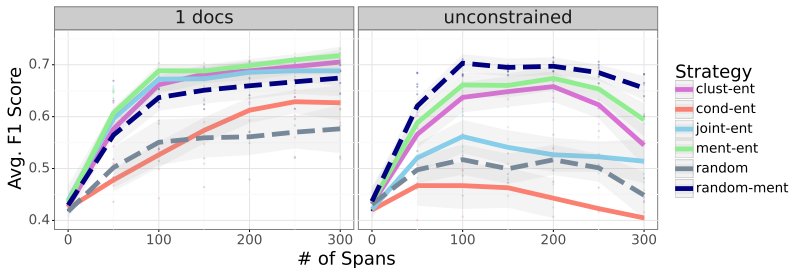
Varying Number of Documents Read

- Test Avg. F1 on PreCo of each strategy
- On each cycle, sample fifty spans from either one document or across many documents



Varying Number of Documents Read

- Test Avg. F1 on PreCo of each strategy
- On each cycle, sample fifty spans from either one document or across many documents



Model training is unstable with unconstrained sampling

Reading and Labeling for Humans

Three users label spans sampled from PreCo

Text

Two new studies have investigated why fewer women , compared to men , study and work in the so-called STEM subjects in the United States : science , technology , engineering and mathematics . The American Association of University Women (AAUW) examined existing research . Its report `` Why So Few ? '' suggested ways to interest more girls and women in the STEM fields . The researchers found that cultural and environmental factors make a difference . Researcher Christianne Corbett says more boys than girls score very high on math

Active query: Researcher Christianne Corbett Answer: The researchers

Queries: the STEM fields The researchers Researcher Christianne Corbett Iceland

Overlapping Candidates

(1) The researchers

(n)o previous mention

(q)uery is not entity

Reading and Labeling for Humans

Three users label spans sampled from PreCo

Text

Two new studies have investigated why fewer women , compared to men , study and work in the so-called STEM subjects in the United States : science , technology , engineering and mathematics . The American Association of University Women (AAUW) examined existing research . Its report `` Why So Few ? '' suggested ways to interest more girls and women in the STEM fields . The researchers found that cultural and environmental factors make a difference . Researcher Christianne Corbett says more boys than girls score very high on math

Active query: Researcher Christianne Corbett Answer: The researchers

Queries: the STEM fields The researchers Researcher Christianne Corbett Iceland

Overlapping Candidates

(1) The researchers

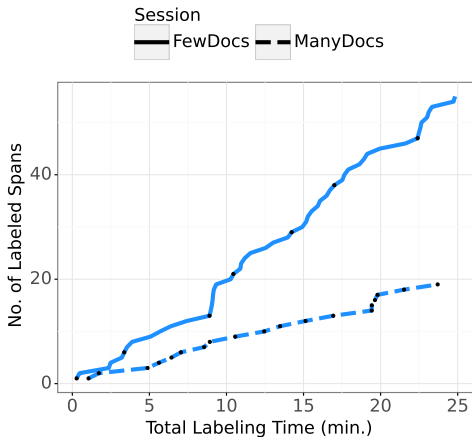
(n)o previous mention

(q)uery is not entity

Users complete two twenty-five minute sessions:

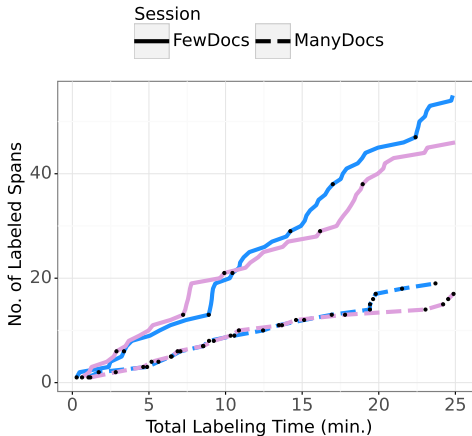
1. **FewDocs:** Read fewer docs. and label multiple spans per doc.
2. **ManyDocs:** Read more docs. and label one span per doc.

Labeling Throughput At Least Doubles in FewDocs



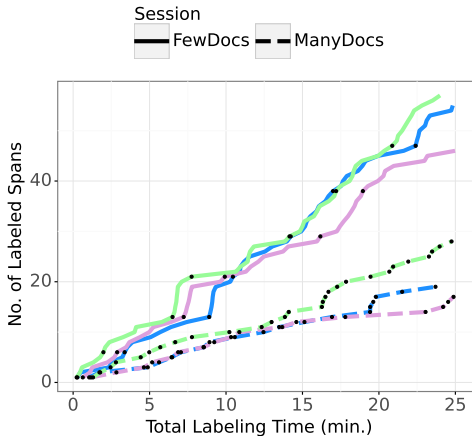
- Each color indicates one of three users and the linetype designates the session
- Black dots mark the first span labeled in a different document

Labeling Throughput At Least Doubles in FewDocs



- Each color indicates one of three users and the linetype designates the session
- Black dots mark the first span labeled in a different document

Labeling Throughput At Least Doubles in FewDocs



- Each color indicates one of three users and the linetype designates the session
- Black dots mark the first span labeled in a different document

Thanks

Any Questions?

myuan@cs.umd.edu



DEPARTMENT OF
COMPUTER SCIENCE

References I

- Xilun Chen, Yu Sun, Ben Athiwaratkun, Claire Cardie, and Kilian Weinberger. Adversarial deep averaging networks for cross-lingual sentiment classification. 6:557–570, 2018. doi: 10.1162/tacI_a_00039.
- Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S Weld, Luke Zettlemoyer, and Omer Levy. SpanBERT: Improving pre-training by representing and predicting spans. 8:64–77, 2020. doi: 10.1162/tacI_a_00300.
- Kenton Lee, Luheng He, and Luke Zettlemoyer. Higher-order coreference resolution with coarse-to-fine inference. 2018. doi: 10.18653/v1/N18-2108.
- Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Hwee Tou Ng, Anders Björkelund, Olga Uryupina, Yuchen Zhang, and Zhi Zhong. Towards robust linguistic analysis using OntoNotes. 2013. URL <https://www.aclweb.org/anthology/W13-3516>.
- Patrick Xia and Benjamin Van Durme. Moving on from OntoNotes: Coreference resolution model transfer. 2021. doi: 10.18653/v1/2021.emnlp-main.425.