# An Annotated Corpus for Machine Reading of Instructions in Wet Lab Protocols

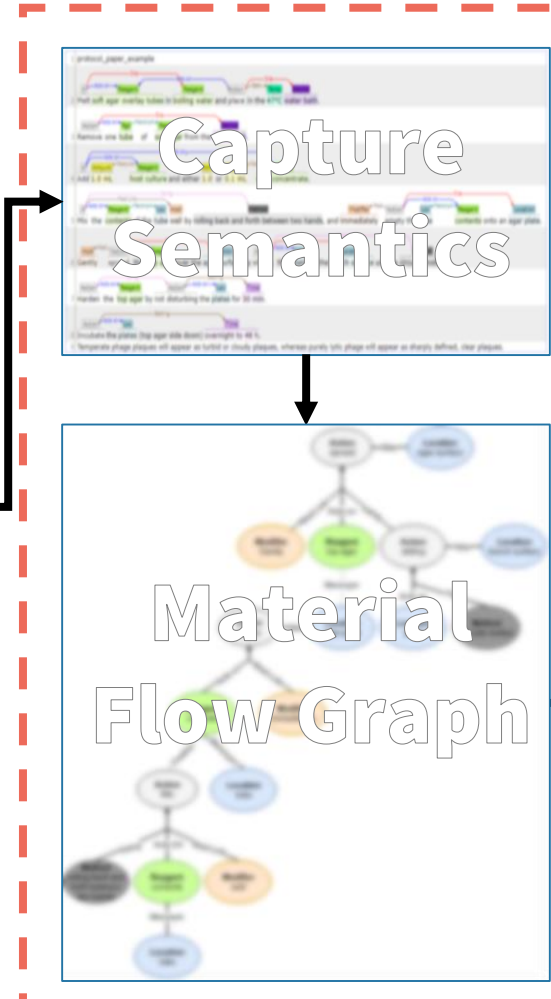Chaitanya Kulkarni, Wei Xu, Alan Ritter, Raghu Machiraju

The Ohio State University
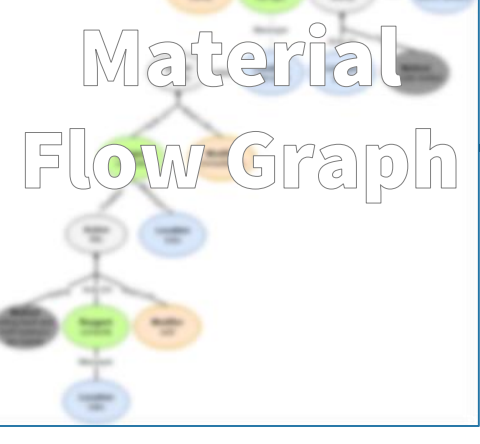
## Introduction

Cumbersome biological experiments necessitates automation to reduce human error and make science reproducible



**Isolation of temperate phages by plaque agar overlay**
1. Melt soft agar overlay tubes in boiling water and place in the 47° C water bath.
2. Remove one tube of soft agar from the water bath.
3. Add 1.0 mL host culture and either 1.0 or 0.1 mL viral concentrate.
4. Mix the contents of the tube well by rolling back and forth between two hands, and immediately empty the tube contents onto an agar plate.
5. Sit RT for 5 min.
6. Gently spread the top agar over the agar surface by sliding the plate on the bench surface using a circular motion.
7. Harden the top agar by not disturbing the plates for 30 min.
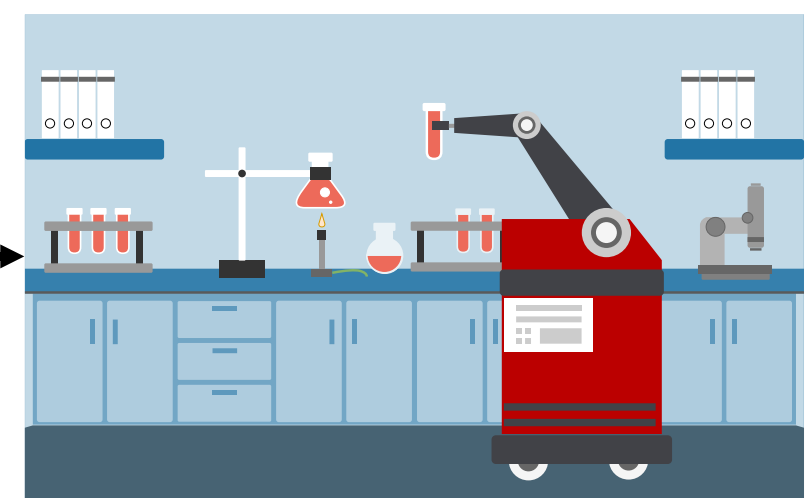8. Incubate the plates (top agar side down) overnight to 48 h.

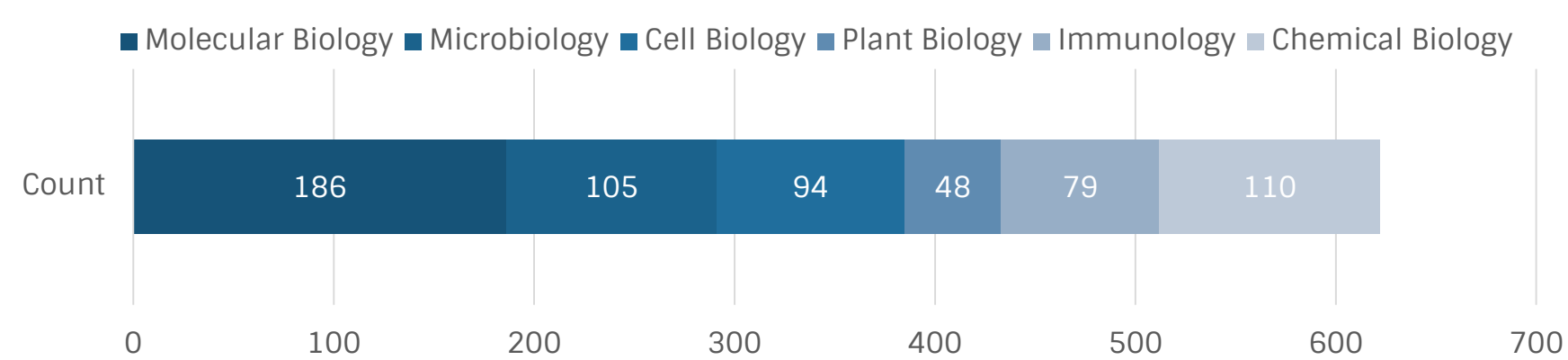The goal is to fill this gap between Natural language instructions and full robotic automation.

We take the first steps towards this goal by,
- Introducing canonical semantic representations understood by experts and non-experts and create a comprehensive corpora, WLP
- Demonstrating utility of corpora by developing machine learning approaches for semantic parsing of instructions
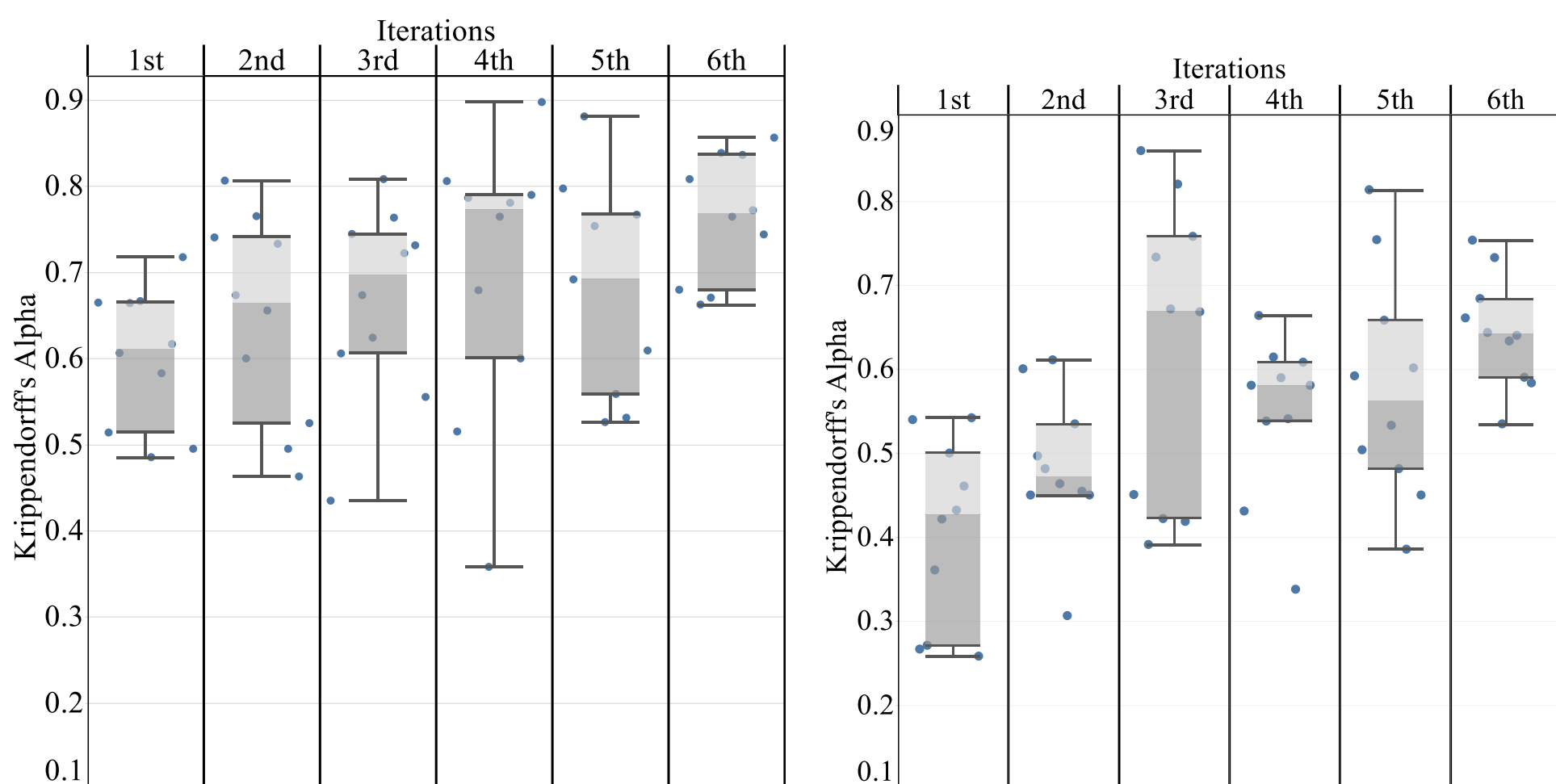
Web Demo

## Corpus Statistics



|  | Total | per Protocol | per Sentence |
|---|---|---|---|
| # of sentences | 13679 | 21.99 | – |
| # of words | 177770 | 285.80 | 12.99 |
| # of entities | 43236 | 69.51 | 3.16 |
| # of relations | 42425 | 68.21 | 3.10 |
| # of actions | 17485 | 28.11 | 1.28 |

Our corpus consist of 622 protocols annotated by a team of 10 annotators

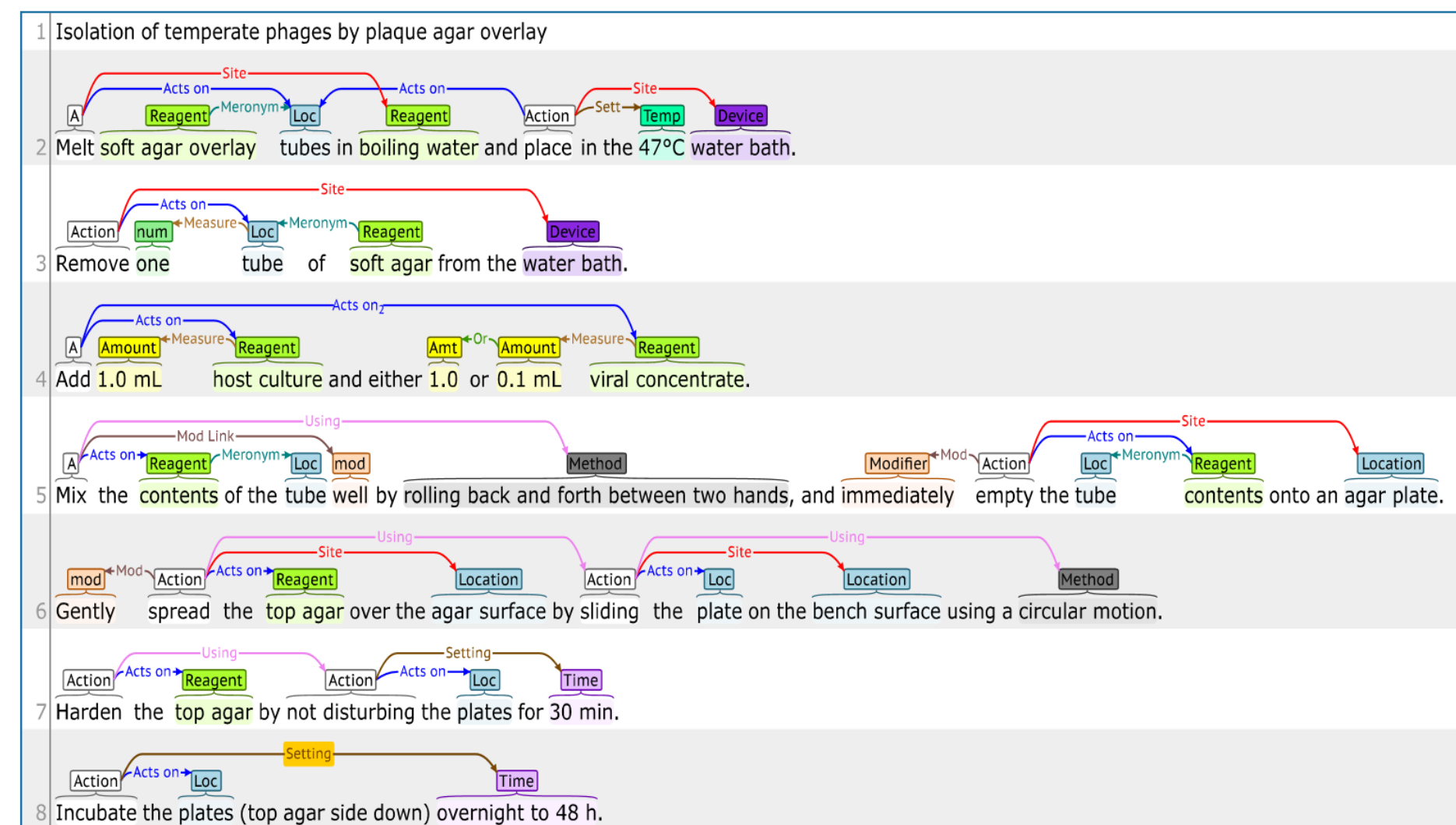## Inter-Annotator Agreement



Action + Entities          Relations

Every iteration consisted of 10 protocols, annotated by 4 coders.

Inter-annotator agreement improves over iterations as changes were made in the annotation guidelines.

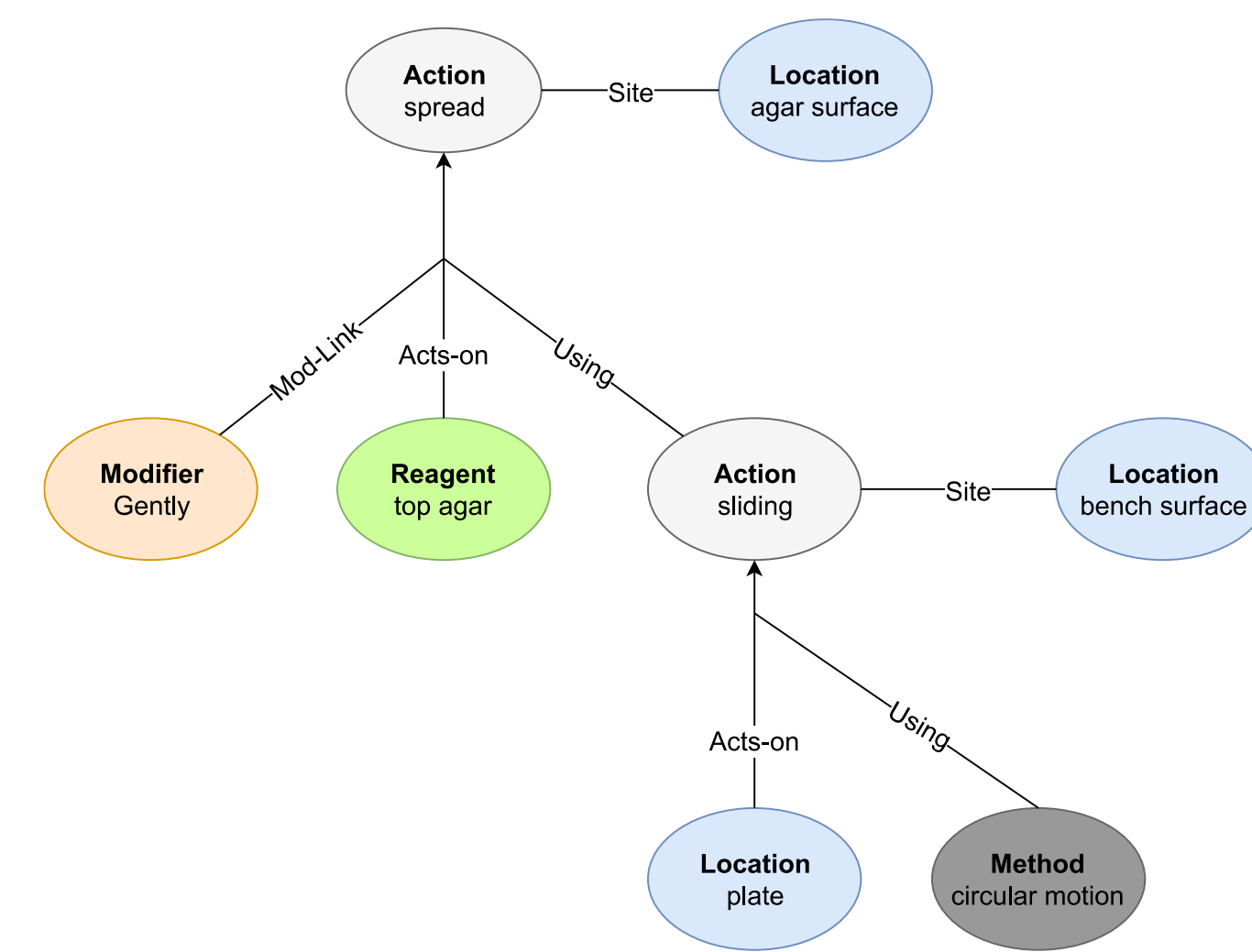| Annotators | Entities+Actions | Relations |
|---|---|---|
| Biologist-Linguist | 0.7600 | 0.6084 |
| Biologist-Other | 0.7621 | 0.6619 |
| Linguist-Other | 0.7574 | 0.6753 |
| all 4 coders | 0.7599 | 0.6625 |

Similar inter annotator agreement between annotators with varying backgrounds demonstrates the ease of comprehension from experts and non-experts alike.
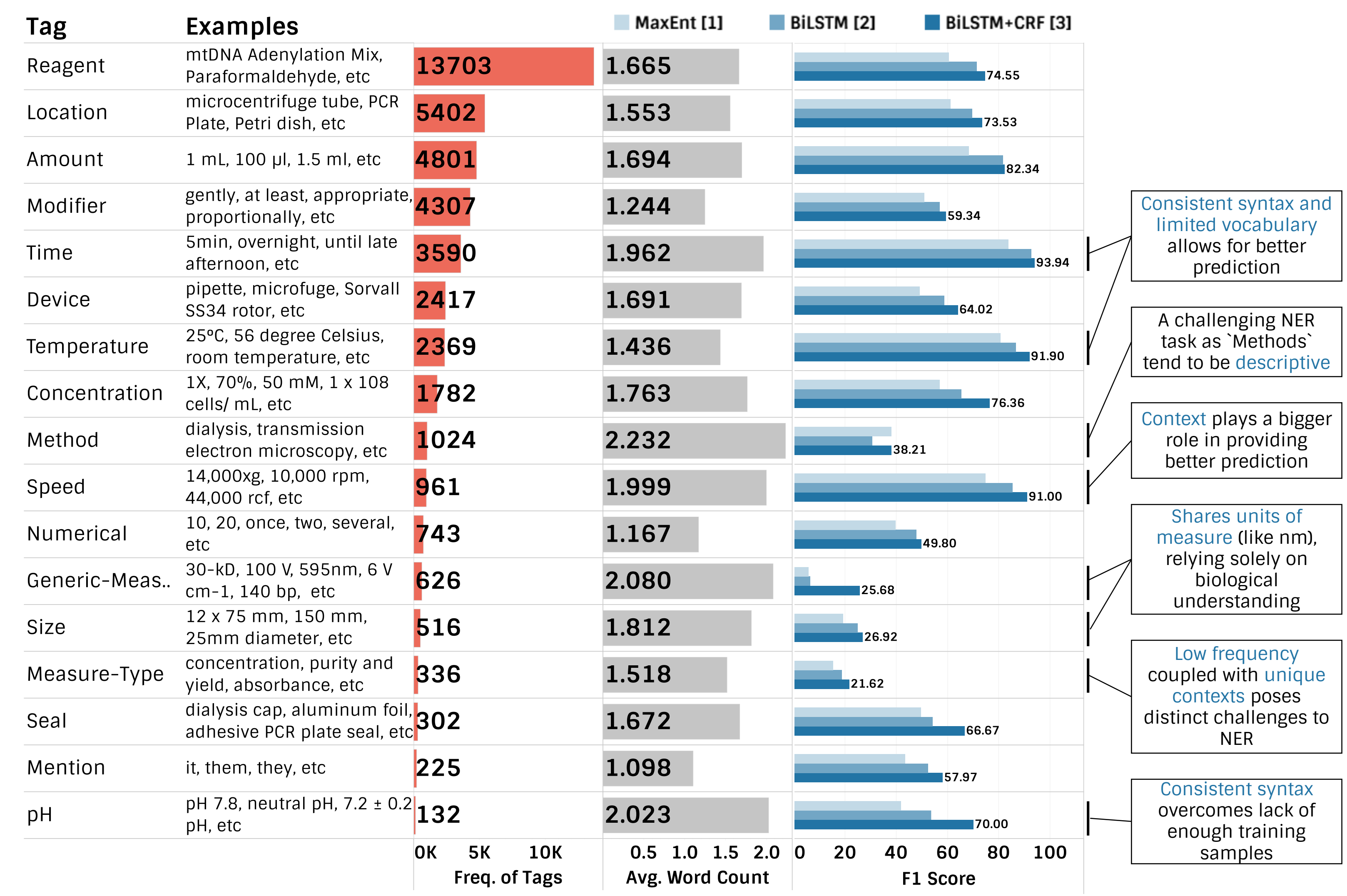
## Wet Lab Protocol Corpus



Wet lab protocols are sequence of steps consisting of:
- **Imperative Sentences**: instructing actions
- **Declarative Sentences**: describing result of a previous action
- **Notes**: general guidelines and/or warnings

WLP corpus constructs canonical representation for a wet lab protocol – an action graph directly derived from the annotations.
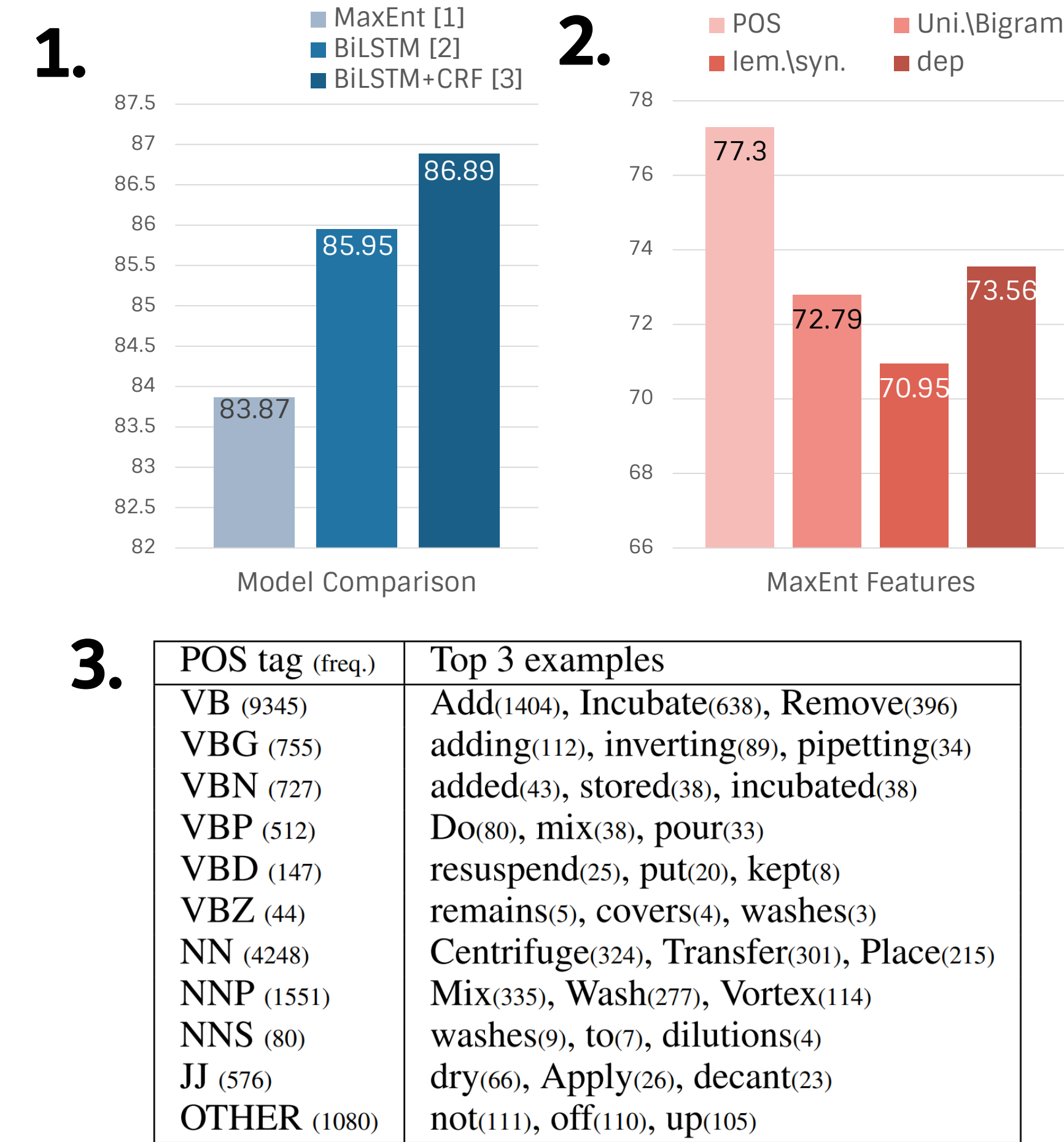


## Entity Extraction

| Tag | Examples | MaxEnt [1] | BiLSTM [2] | BiLSTM+CRF [3] |
|---|---|---|---|---|
| Reagent | mtDNA Adenylation Mix, Paraformaldehyde, etc | 13703 | 1.665 | 74.55 |
| Location | microcentrifuge tube, PCR Plate, Petri dish, etc | 5402 | 1.553 | 73.53 |
| Amount | 1 mL, 100 µl, 1.5 ml, etc | 4801 | 1.694 | 82.34 |
| Modifier | gently, at least, appropriate, proportionally, etc | 4307 | 1.244 | 59.34 |
| Time | 5min, overnight, until late afternoon, etc | 3590 | 1.962 | 93.94 |
| Device | pipette, microfuge, Sorvall SS34 rotor, etc | 2417 | 1.691 | 64.02 |
| Temperature | 25°C, 56 degree Celsius, room temperature, etc | 2369 | 1.436 | 91.90 |
| Concentration | 1X, 70%, 50 mM, 1 x 108 cells/ mL, etc | 1782 | 1.763 | 76.36 |
| Method | dialysis, transmission electron microscopy, etc | 1024 | 2.232 | 38.21 |
| Speed | 14,000xg, 10,000 rpm, 44,000 rcf, etc | 961 | 1.999 | 91.00 |
| Numerical | 10, 20, once, two, several, etc | 743 | 1.167 | 49.80 |
| Generic-Meas.. | 30-kD, 100 V, 595nm, 6 V cm-1, 140 bp, etc | 626 | 2.080 | 25.68 |
| Size | 12 x 75 mm, 150 mm, 25mm diameter, etc | 516 | 1.812 | 26.92 |
| Measure-Type | concentration, purity and yield, absorbance, etc | 336 | 1.518 | 21.62 |
| Seal | dialysis cap, aluminum foil, adhesive PCR plate seal, etc | 302 | 1.672 | 66.67 |
| Mention | it, them, they, etc | 225 | 1.098 | 57.97 |
| pH | pH 7.8, neutral pH, 7.2 ± 0.2 pH, etc | 132 | 2.023 | 70.00 |

Freq. of Tags / Avg. Word Count / F1 Score

- Consistent syntax and limited vocabulary allows for better prediction
- A challenging NER task as 'Methods' tend to be descriptive
- Context plays a bigger role in providing better prediction
- Shares units of measure (like nm), relying solely on biological understanding
- Low frequency coupled with unique contexts poses distinct challenges to NER
- Consistent syntax overcomes lack of enough training samples

## References

1. Andrew Borthwick and Ralph Grishman. 1999. A maximum entropy approach to named entity recognition. *Ph. D. Thesis, Dept. of Computer Science, New York University.*
2. Alan Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. 2013. Speech recognition with deep recurrent neural networks. *In Proceedings of the 2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 6645–6649.
3. Xuezhe Ma and Eduard Hovy. 2016. End-to-end se-quence labeling via bi-directional lstm-cnns-crf. *In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL).*
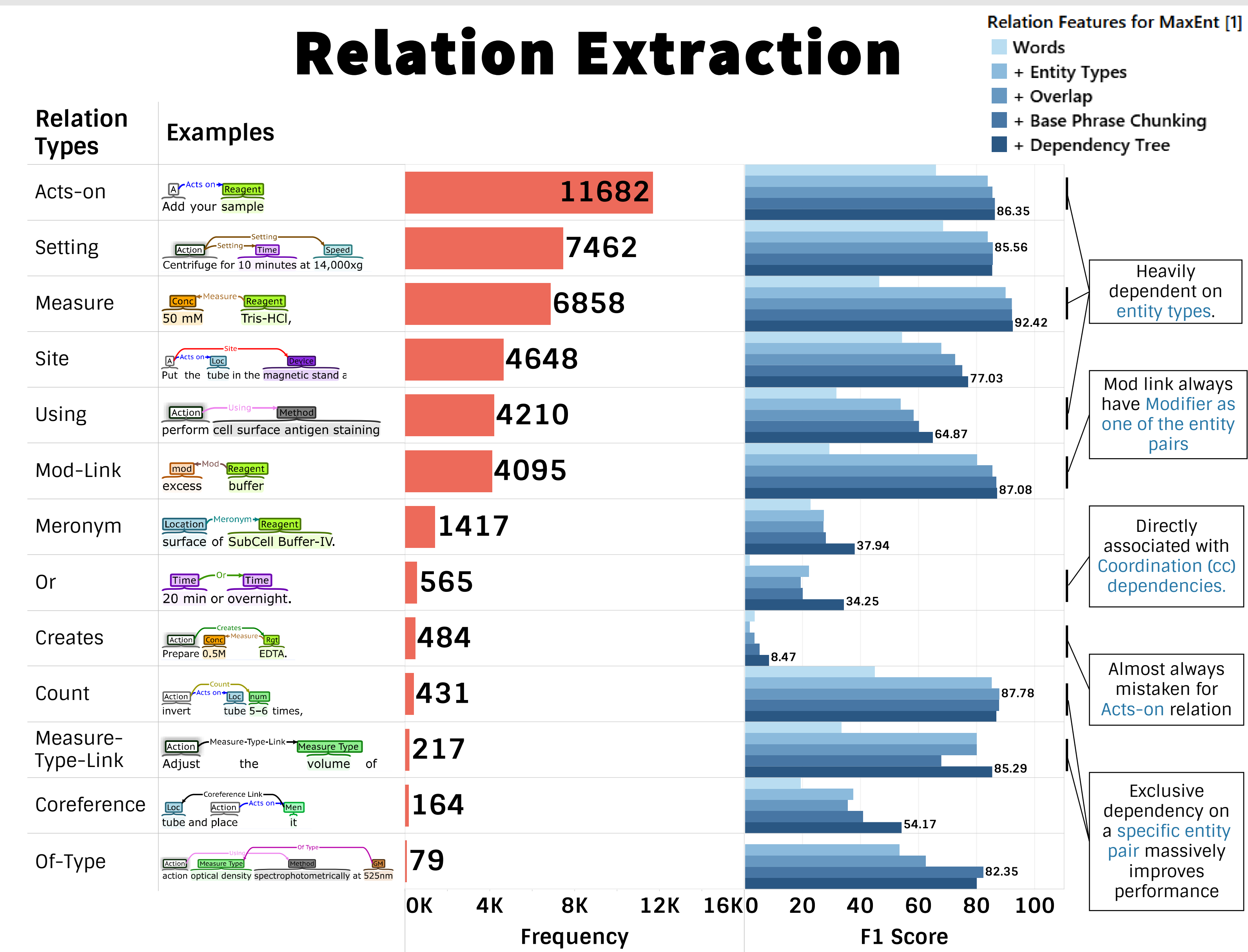
## Action Extraction



1. Evaluate corpora by classifying actions using the best maximum entropy model and 2 neural models.

2. Parts of speech was the most effective in capturing action words

3. Majority of the action verbs fall under **VBs (60.48%)** or **NN (30.84%)** using GENIA POS tagger. A small percentage are misclassified under **OTHER (5.66%)** and **JJ (3.02%)**.

| POS tag (freq.) | Top 3 examples |
|---|---|
| VB (9345) | Add(1404), Incubate(638), Remove(396) |
| VBG (755) | adding(112), inverting(89), pipetting(34) |
| VBN (727) | added(43), stored(38), incubated(38) |
| VBP (512) | Do(80), mix(38), pour(33) |
| VBD (147) | resuspend(25), put(20), kept(8) |
| VBZ (44) | remains(5), covers(4), washes(3) |
| NN (4248) | Centrifuge(324), Transfer(301), Place(215) |
| NNP (1551) | Mix(335), Wash(277), Vortex(114) |
| NNS (80) | washes(9), to(7), dilutions(4) |
| JJ (576) | dry(66), Apply(26), decant(23) |
| OTHER (1080) | not(111), off(110), up(105) |

## Relation Extraction



Relation Features for MaxEnt [1]: Words, + Entity Types, + Overlap, + Base Phrase Chunking, + Dependency Tree

| Relation Types | Examples | Frequency | F1 Score |
|---|---|---|---|
| Acts-on | Add your sample | 11682 | 86.35 |
| Setting | Centrifuge for 10 minutes at 14,000xg | 7462 | 85.56 |
| Measure | 50 mM Tris-HCl, | 6858 | 92.42 |
| Site | Put the tube in the magnetic stand a | 4648 | 77.03 |
| Using | perform cell surface antigen staining | 4210 | 64.87 |
| Mod-Link | excess buffer | 4095 | 87.08 |
| Meronym | surface of SubCell Buffer-IV. | 1417 | 37.94 |
| Or | 20 min or overnight. | 565 | 34.25 |
| Creates | Prepare 0.5M EDTA. | 484 | 8.47 |
| Count | invert tube 5–6 times, | 431 | 87.78 |
| Measure-Type-Link | Adjust the volume of | 217 | 85.29 |
| Coreference | tube and place it | 164 | 54.17 |
| Of-Type | action optical density spectrophotometrically at 525nm | 79 | 82.35 |

- Heavily dependent on entity types.
- Mod link always have Modifier as one of the entity pairs
- Directly associated with Coordination (cc) dependencies.
- Almost always mistaken for Acts-on relation
- Exclusive dependency on a specific entity pair massively improves performance

## Conclusion + Future Work

We present a corpus with accessible semantic representation. Given the varying emphasis on morphology and context, every named entity and relation in this representation poses unique challenges to semantic parsing.

In addition to implementing methods that address these challenges, we plan to extend the corpus by inter connecting all the sentences to build a more complete protocol representation.