

# Deep Just-In-Time Inconsistency Detection Between Comments and Source Code

Sheena Panthaplackel, Junyi Jessy Li, Milos Gligoric,  
Raymond J. Mooney

# Source Code Comments

**Document functionality,  
usage, implementation, error  
cases, ...**

```
/**Computes the highest value from the list of scores.*/  
public int getBestScore() {  
    return Collections.max(scores);  
}
```

# Source Code Comments

**When developers make code changes, they often fail to update comments accordingly**

```
/**Computes the highest value from the list of scores.*/  
public int getBestScore() {  
    return Collections.max(scores);  
}
```

- return Collections.max(scores);  
+ return Collections.min(scores);



Leads to time-wasting confusion  
and vulnerability to bugs

```
/**Computes the highest value from the list of scores.*/  
public int getBestScore() {  
    return Collections.min(scores);  
}
```

**Automatically detect inconsistent comments**

# Detecting Inconsistent Comments

## Just-In-Time Inconsistency Detection

```
/**Computes the highest value from the list of scores.*/  
public int getBestScore() {  
    return Collections.max(scores);  
}
```

- return Collections.max(scores);  
+ return Collections.min(scores);

```
/**Computes the highest value from the list of scores.*/  
public int getBestScore() {  
    return Collections.max(scores);  
}
```

- public int getBestScore() {  
+ public double getBestScore() {

## Post Hoc Inconsistency Detection

```
/**Computes the highest value from the list of scores.*/  
public int getBestScore() {  
    return Collections.min(scores);  
}
```

Inconsistent 

```
/**Computes the highest value from the list of scores.*/  
public double getBestScore() {  
    return Collections.max(scores);  
}
```

Consistent 

# Outline

- Task
- Architecture
- Data
- Intrinsic Evaluation
- Extrinsic Evaluation
- Summary

# Outline

- Task
- Architecture
- Data
- Intrinsic Evaluation
- Extrinsic Evaluation
- Summary

# Task

Determine whether a **comment (C)** is inconsistent with a **method (M)**

```
/**Computes the highest value from the list of scores.**/
```

```
public int getBestScore() {  
    return Collections.min(scores);  
}
```

```
public int getBestScore() {  
    return Collections.max(scores);  
}
```

Post Hoc: Given **C** and **M**

Just-In-Time: Given **C**, **M**, and **M<sub>old</sub>**

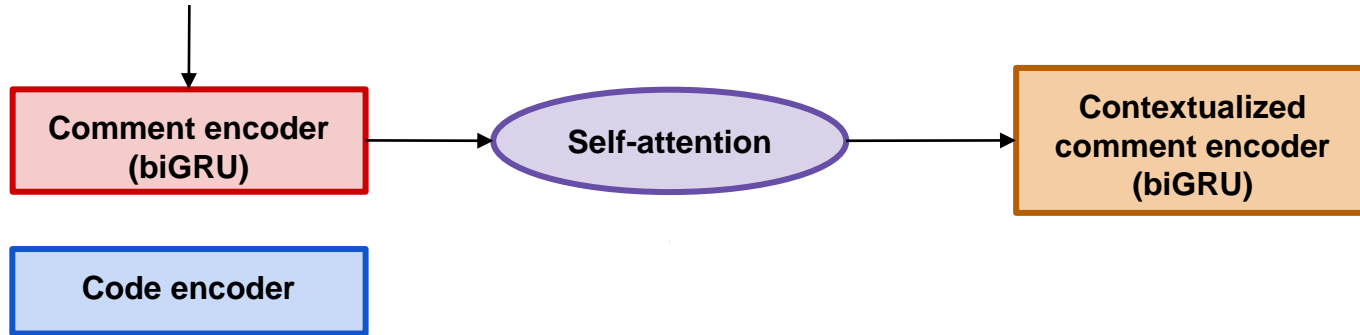
# Outline

- Task
- **Architecture**
- Data
- Intrinsic Evaluation
- Extrinsic Evaluation
- Summary



# Architecture

`/** Computes the highest value from the list of scores. */`



# Code Representations

## Sequence-based

Post hoc (M): M as a sequence of tokens

```
public int getBestScore ( ) { return Collections . min ( scores ) ; }
```

Just-In-Time (M<sub>edit</sub>): Edits between M<sub>old</sub> and M as a sequence of tokens

```
<Keep> public int getBestScore ( ) { return Collections . <KeepEnd>
```

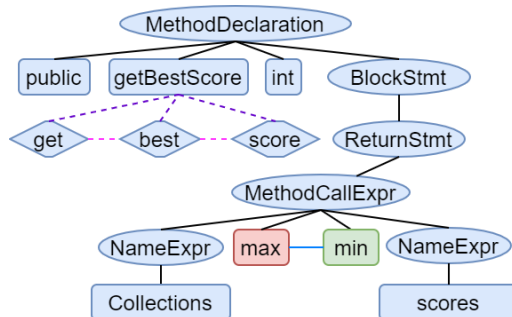
```
<ReplaceOld> max <ReplaceNew> min <ReplaceEnd>
```

```
<Keep> ( scores ) ; } <KeepEnd>
```

## AST-based

Just-In-Time (T<sub>edit</sub>):

Graph representation  
of AST node edits  
between M<sub>old</sub> and M

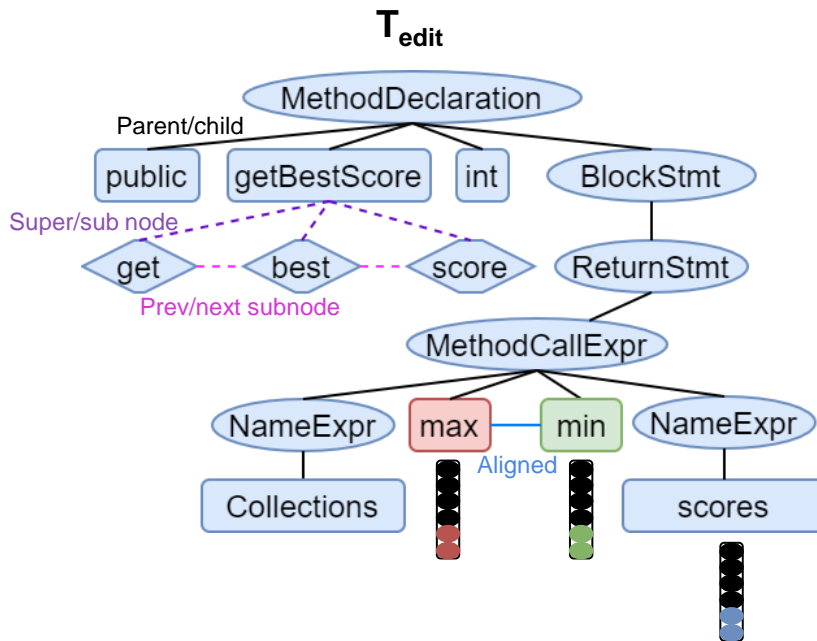


# Encoding Code Representations

- Sequence-based representations ( $\mathbf{M}$ ,  $\mathbf{M}_{\text{edit}}$ ) encoded with a biGRU.
- AST-based representations ( $\mathbf{T}$ ,  $\mathbf{T}_{\text{edit}}$ ) encoded with GGNNs.

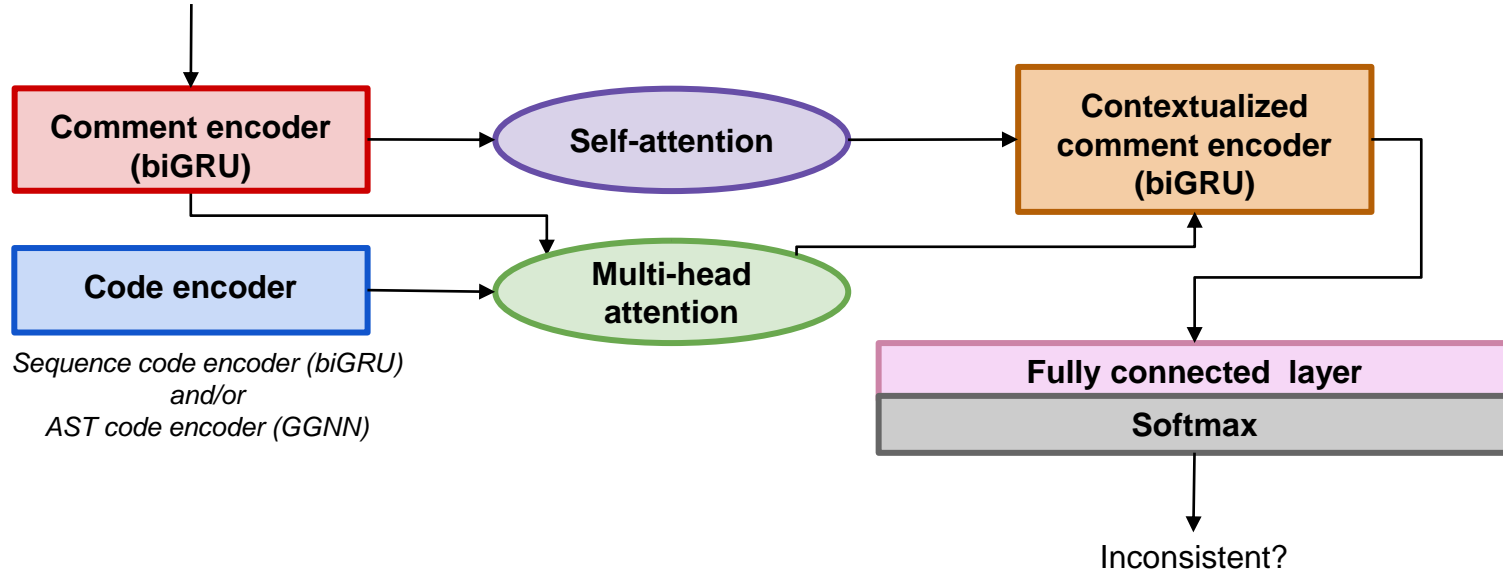
## Gated Graph Neural Network (GGNN)

- Node state representations are updated across 8 steps of message passing
- Message passing is done through different edge types
- Initial state representations consist of a word embedding, concatenated with an edit embedding (for the just-in-time setting)



# Architecture

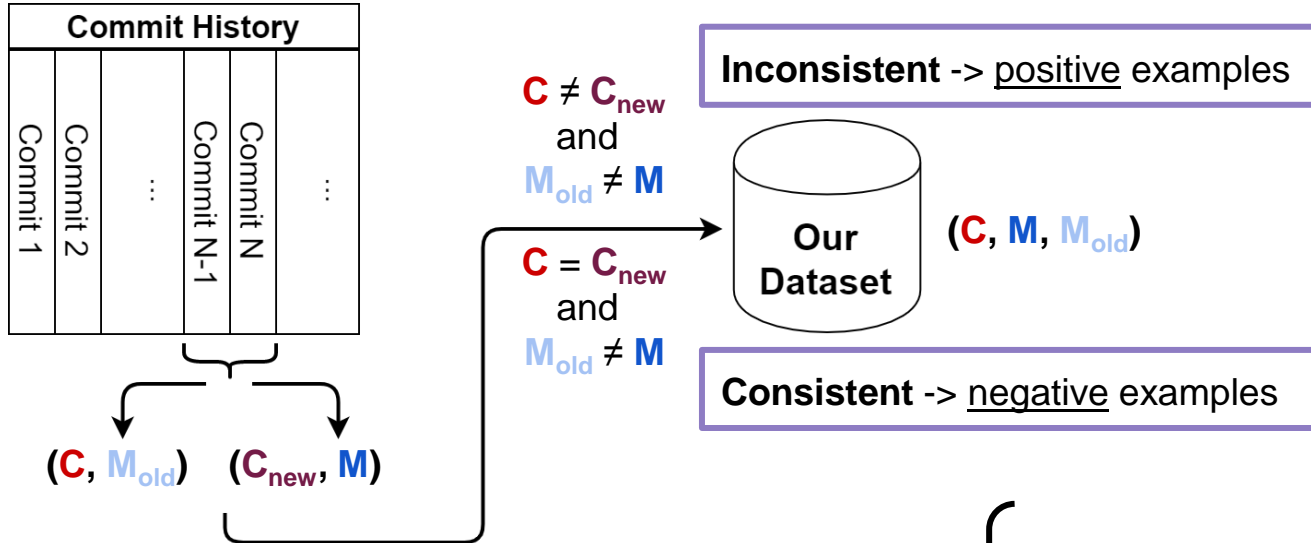
*/\*\* Computes the highest value from the list of scores. \*/*



# Outline

- Task
- Architecture
- **Data**
- Intrinsic Evaluation
- Extrinsic Evaluation
- Summary

# Data Collection



	Train	Valid	Test	Total
Examples	32,988	3,756	3,944	40,688
Projects	829	332	357	1,518

- @return, @param, and summary comments
- Balanced label distribution
- Cleaned 300 examples from test set for more reliable evaluation

# Outline

- Task
- Architecture
- Data
- **Intrinsic Evaluation**
- Extrinsic Evaluation
- Summary

# Intrinsic Evaluation: Results

Cleaned Test Sample

Full Test Set

---

F1

Accuracy

---

F1

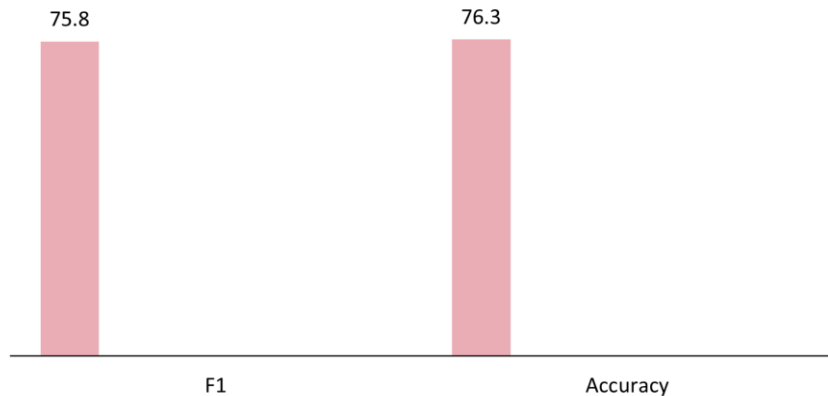
Accuracy



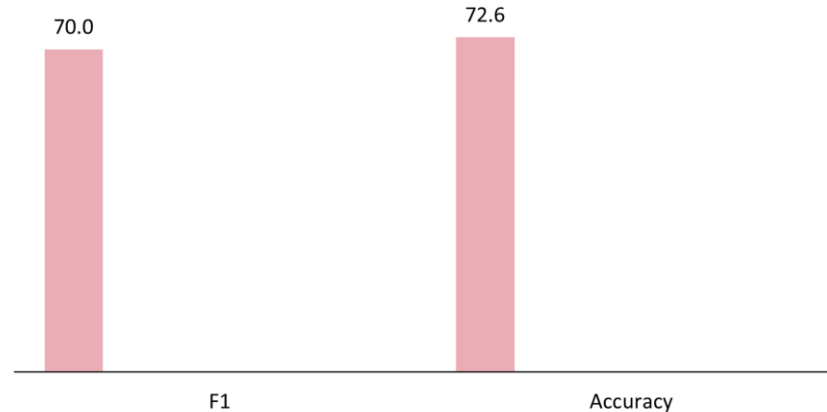
# Intrinsic Evaluation: Results

■ Liu et al. (2018)

Cleaned Test Sample



Full Test Set

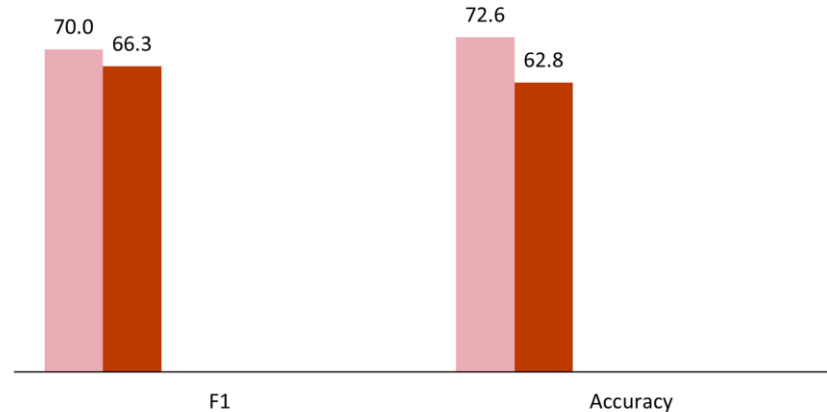
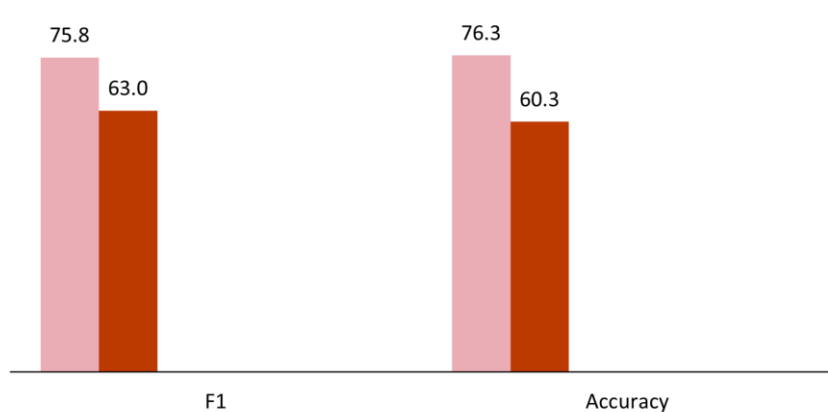


# Intrinsic Evaluation: Results

■ Liu et al. (2018) ■ Post Hoc SEQ

Cleaned Test Sample

Full Test Set

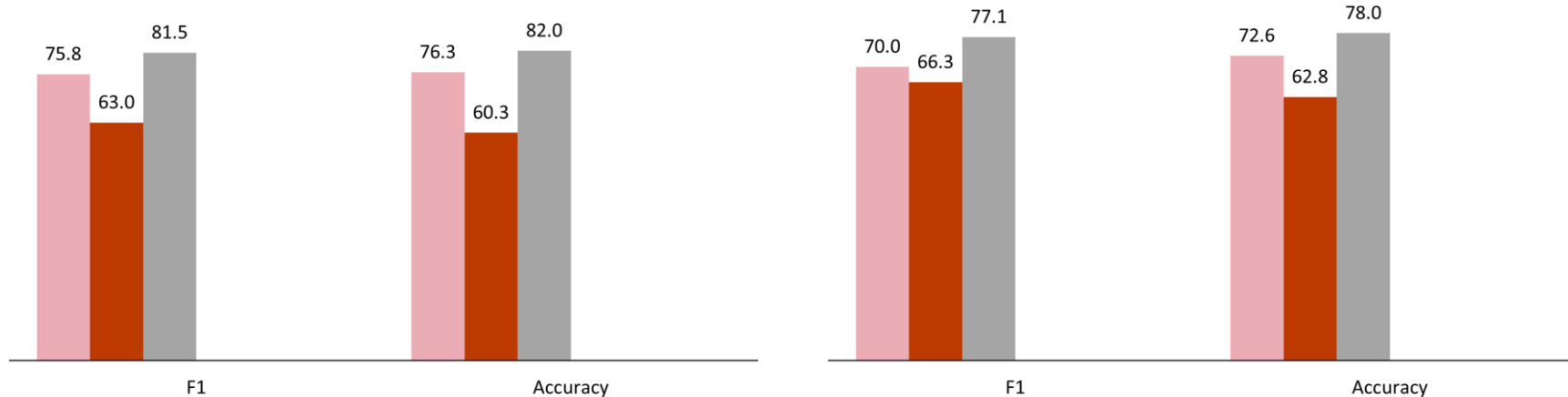


# Intrinsic Evaluation: Results

■ Liu et al. (2018) ■ Post Hoc SEQ ■ Just-In-Time SEQ

Cleaned Test Sample

Full Test Set



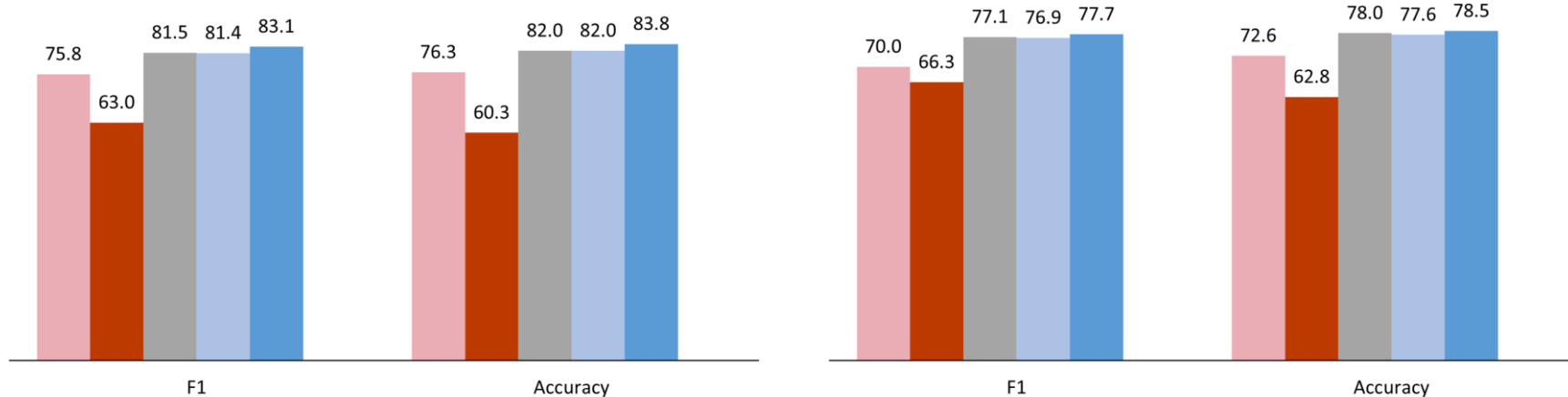
- Our Just-In-Time approach can outperform post hoc and baseline models

# Intrinsic Evaluation: Results

■ Liu et al. (2018)  
 ■ Post Hoc SEQ  
 ■ Just-In-Time SEQ  
 ■ Just-In-Time GRAPH  
 ■ Just-In-Time HYBRID

Cleaned Test Sample

Full Test Set



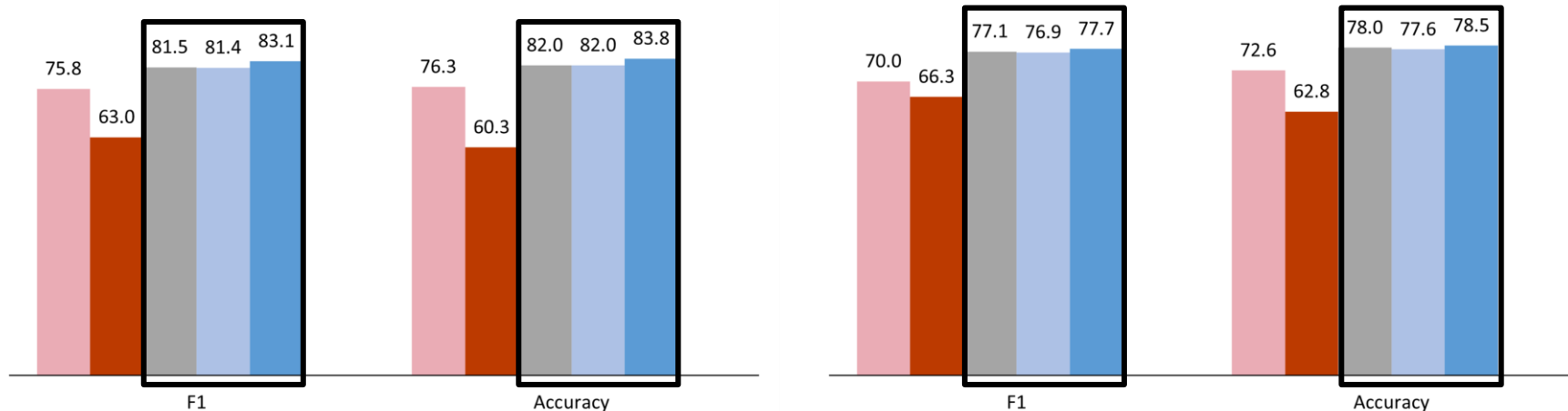
- Our Just-In-Time approach can outperform post hoc and baseline models

# Intrinsic Evaluation: Results

■ Liu et al. (2018)  
 ■ Post Hoc SEQ  
 ■ Just-In-Time SEQ  
 ■ Just-In-Time GRAPH  
 ■ Just-In-Time HYBRID

Cleaned Test Sample

Full Test Set



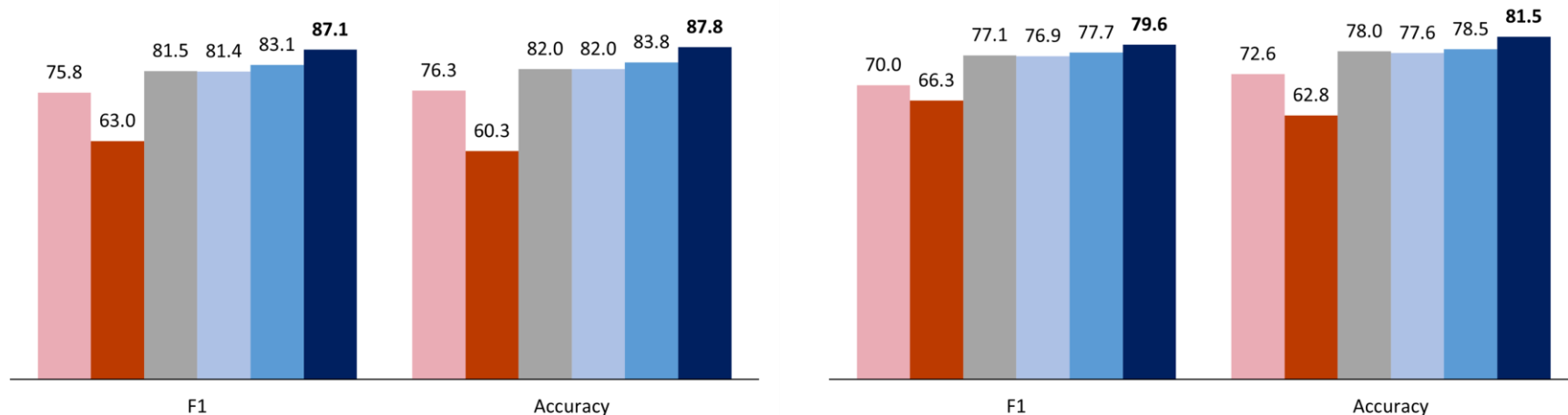
- Our Just-In-Time approach can outperform post hoc and baseline models
- No significant difference between SEQ, GRAPH, and HYBRID approaches

# Intrinsic Evaluation: Results

■ Liu et al. (2018)  
 ■ Post Hoc SEQ  
 ■ Just-In-Time SEQ  
 ■ Just-In-Time GRAPH  
 ■ Just-In-Time HYBRID  
 ■ Just-In-Time HYBRID + features

Cleaned Test Sample

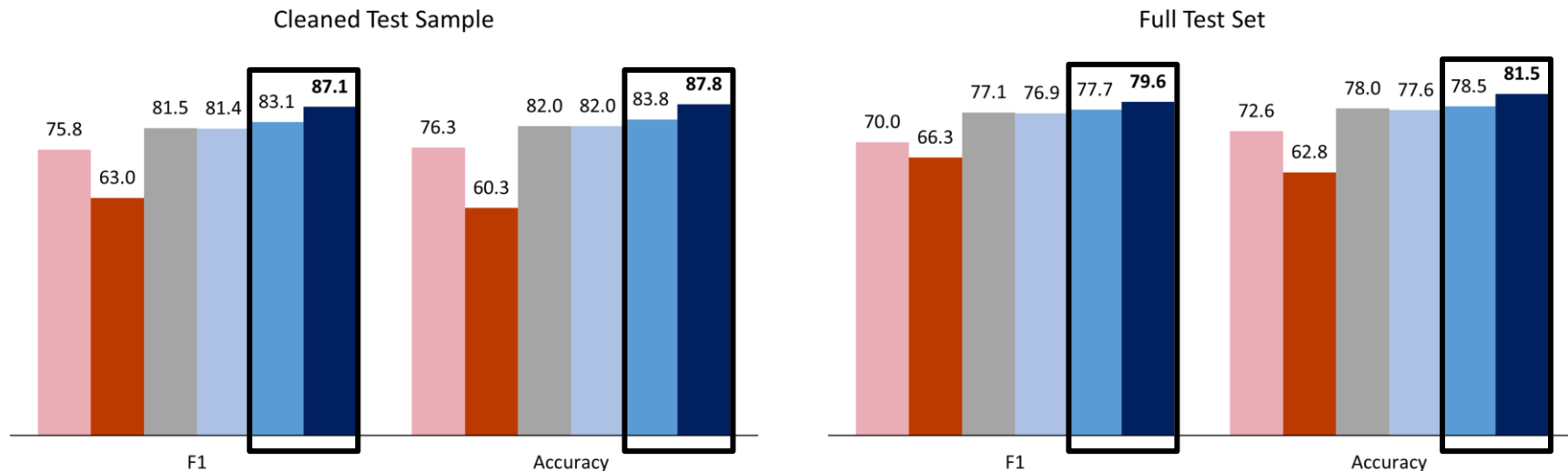
Full Test Set



- Our Just-In-Time approach can outperform post hoc and baseline models
- No significant difference between SEQ, GRAPH, and HYBRID approaches

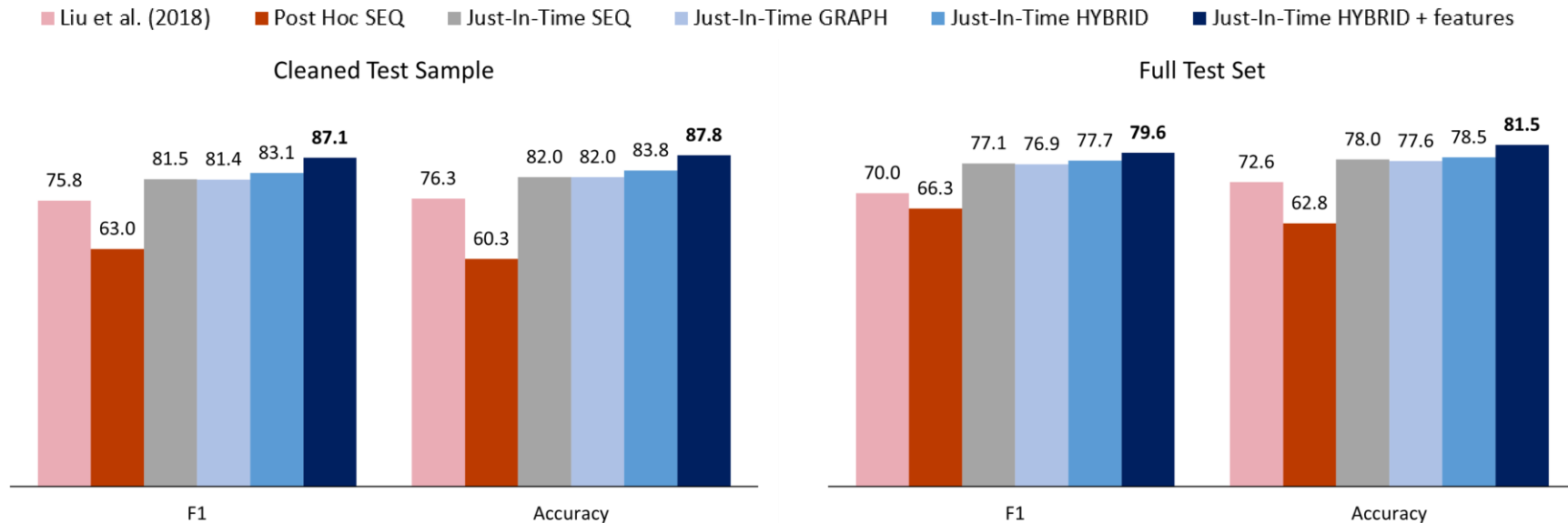
# Intrinsic Evaluation: Results

■ Liu et al. (2018)  
 ■ Post Hoc SEQ  
 ■ Just-In-Time SEQ  
 ■ Just-In-Time GRAPH  
 ■ Just-In-Time HYBRID  
 ■ Just-In-Time HYBRID + features



- Our Just-In-Time approach can outperform post hoc and baseline models
- No significant difference between SEQ, GRAPH, and HYBRID approaches
- Incorporating auxiliary features can further boost performance

# Intrinsic Evaluation: Results



- Our Just-In-Time approach can outperform post hoc and baseline models
- No significant difference between SEQ, GRAPH, and HYBRID approaches
- Incorporating auxiliary features can further boost performance
- Analogous performance between cleaned and full test sets



# Integrating with Update

**On its own, inconsistency detection can only flag comments that developers failed to update.**

```
/**Computes the highest value from the list of scores.*/  
public int getBestScore() {  
    return Collections.max(scores);  
}
```

- return Collections.max(scores);

+ return Collections.min(scores);

**Automatically detect and update inconsistent comments**



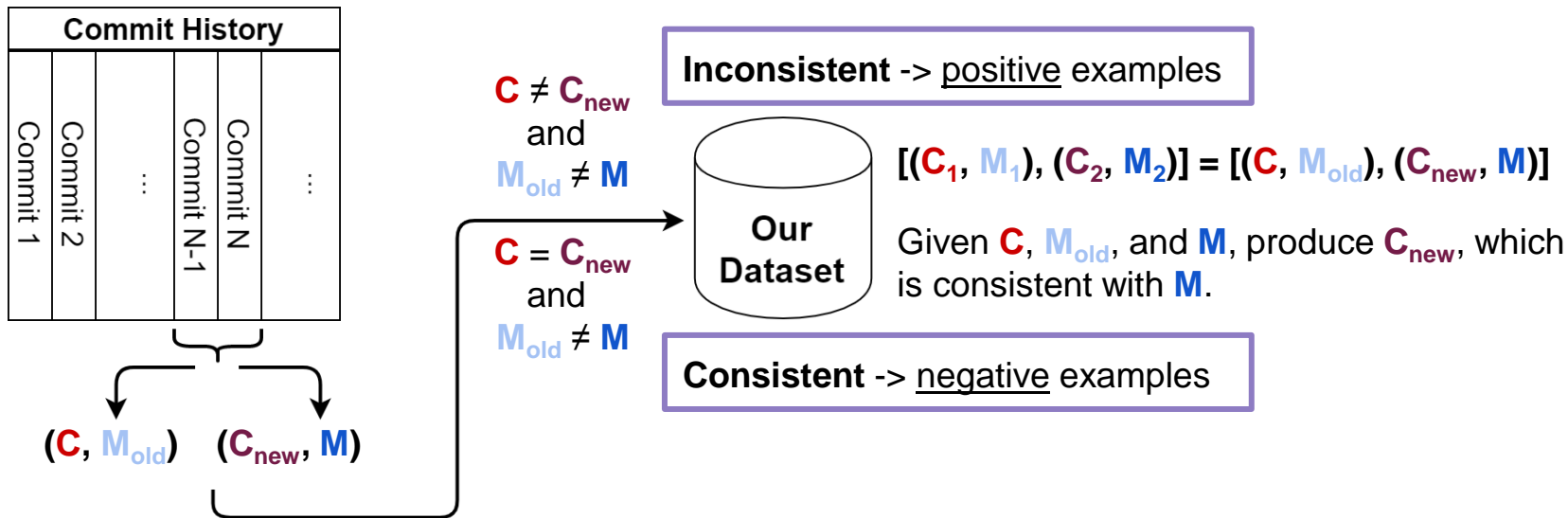
**lowest**

```
/**Computes the highest value from the list of scores.*/  
public int getBestScore() {  
    return Collections.min(scores);  
}
```

# Outline

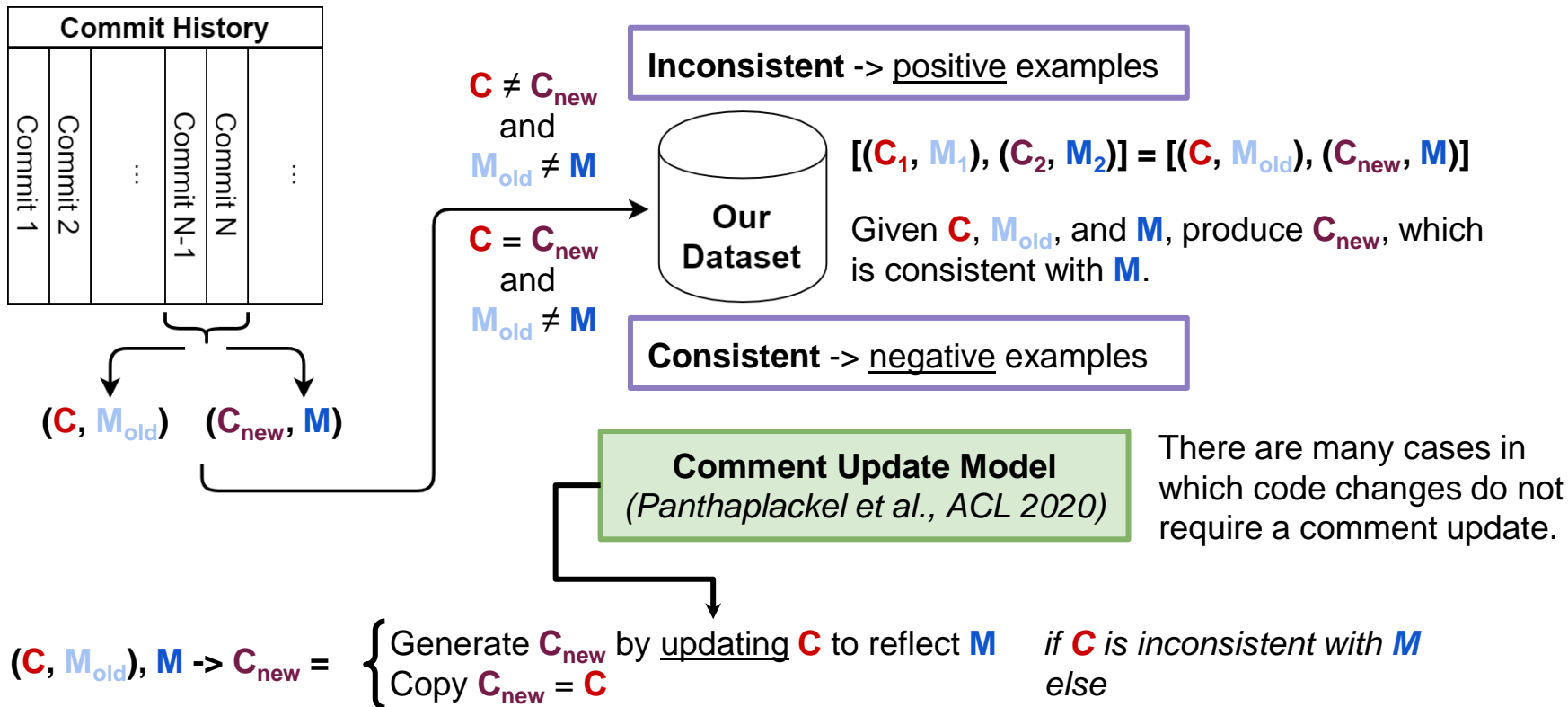
- Task
- Architecture
- Data
- Intrinsic Evaluation
- **Extrinsic Evaluation**
- Summary

# Extrinsic Evaluation: Integrating with Update

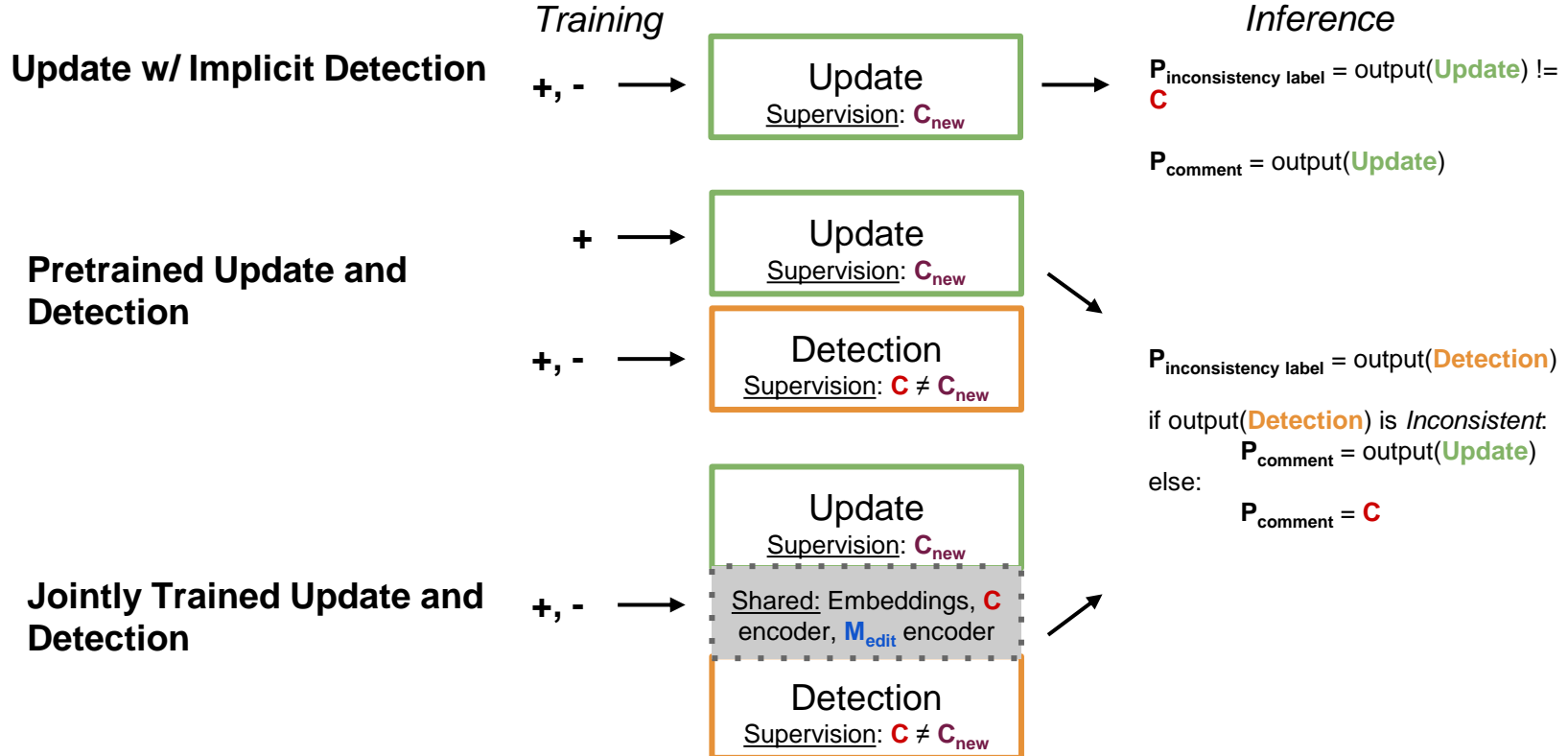


$$(C, M_{old}), M \rightarrow C_{new} = \begin{cases} \text{Generate } C_{new} \text{ by updating } C \text{ to reflect } M & \text{if } C \text{ is inconsistent with } M \\ \text{Copy } C_{new} = C & \text{else} \end{cases}$$

# Extrinsic Evaluation: Integrating with Update



# Extrinsic Evaluation: Integrating with Update

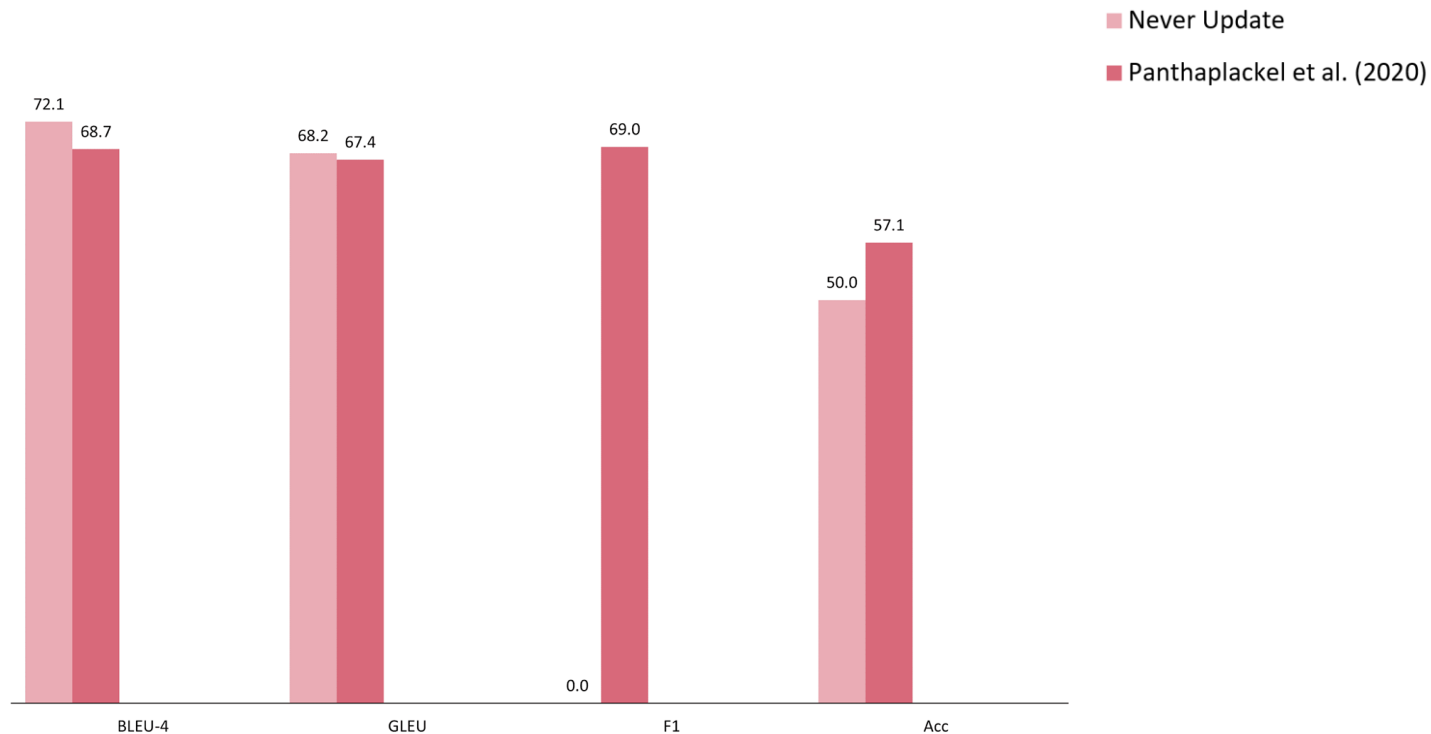


# Extrinsic Evaluation: Results

Cleaned Test Sample

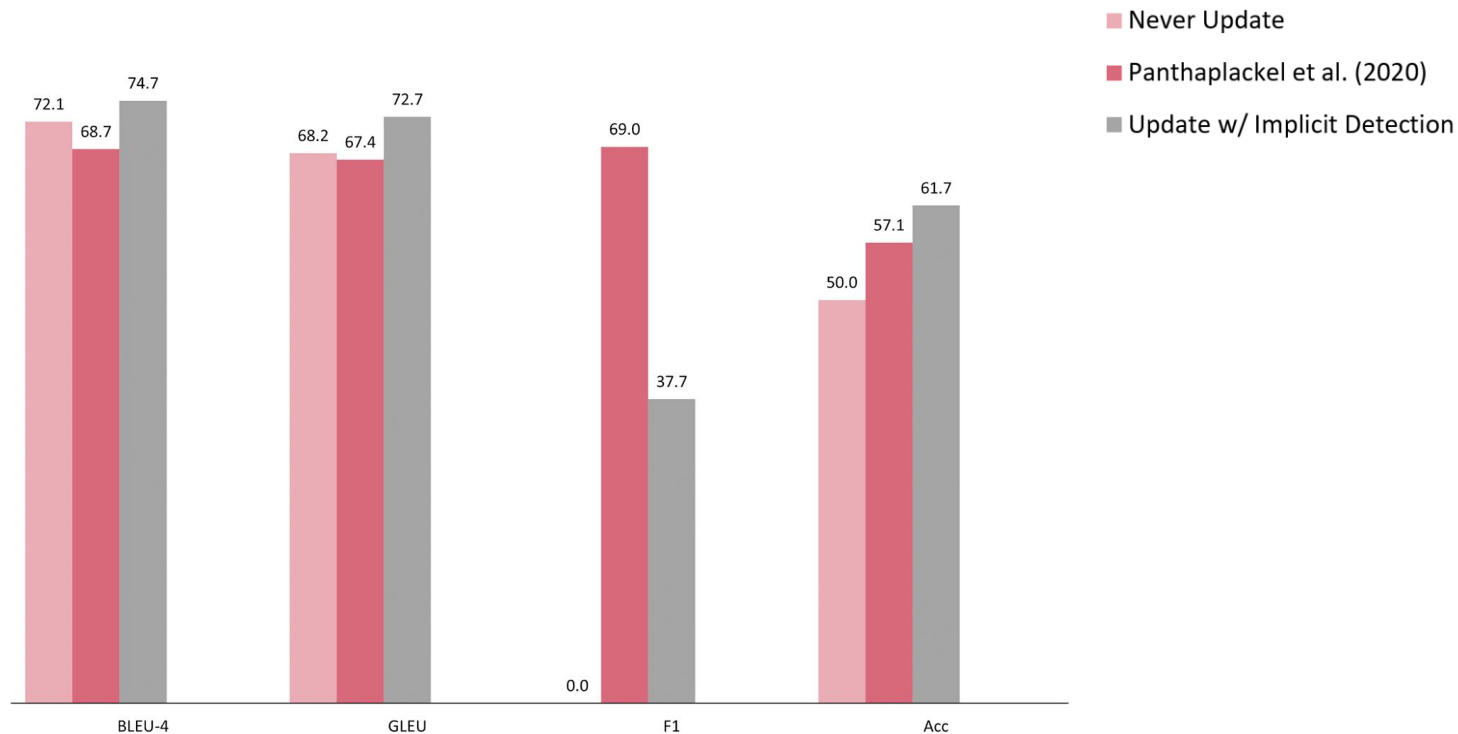
# Extrinsic Evaluation: Results

Cleaned Test Sample



# Extrinsic Evaluation: Results

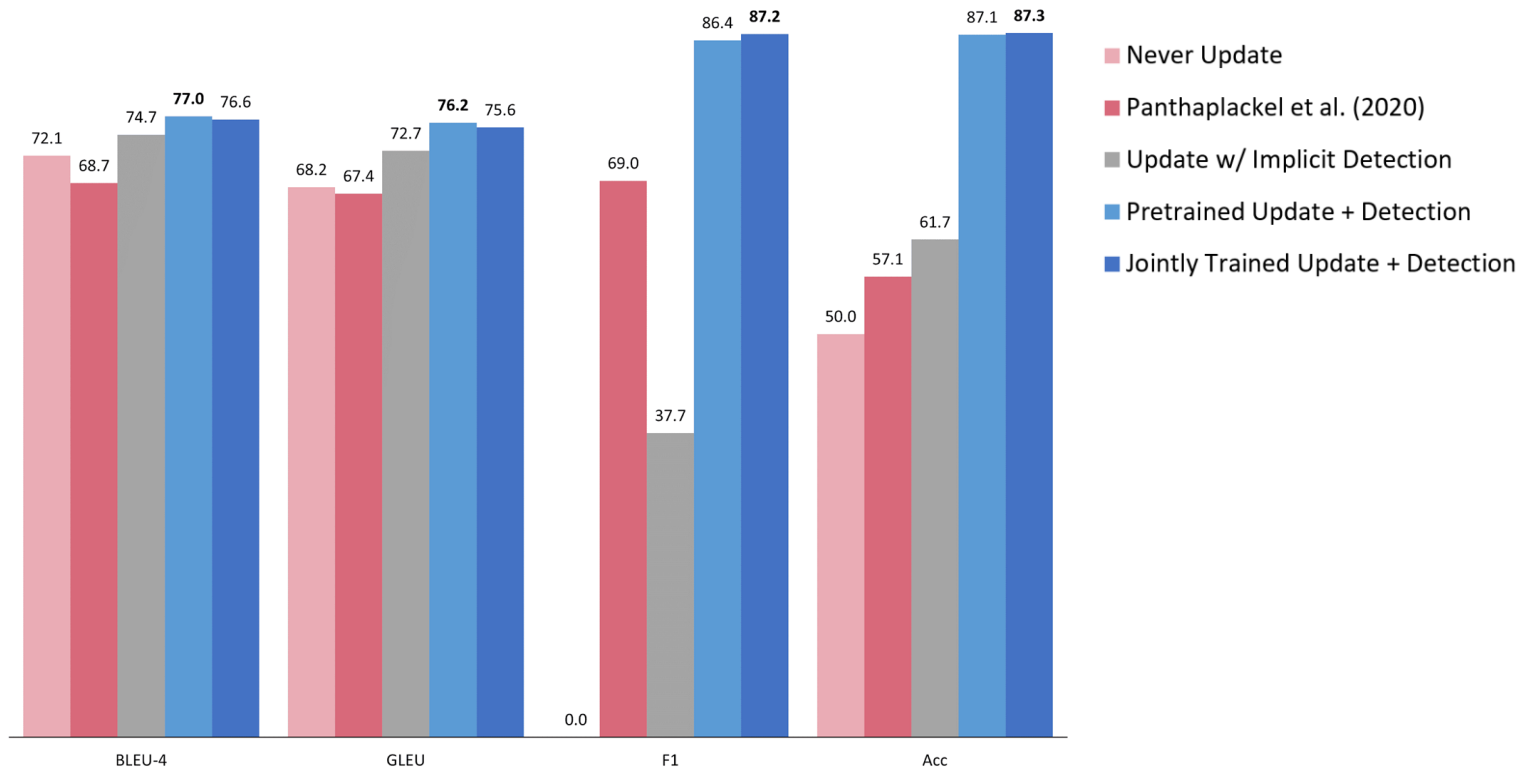
Cleaned Test Sample





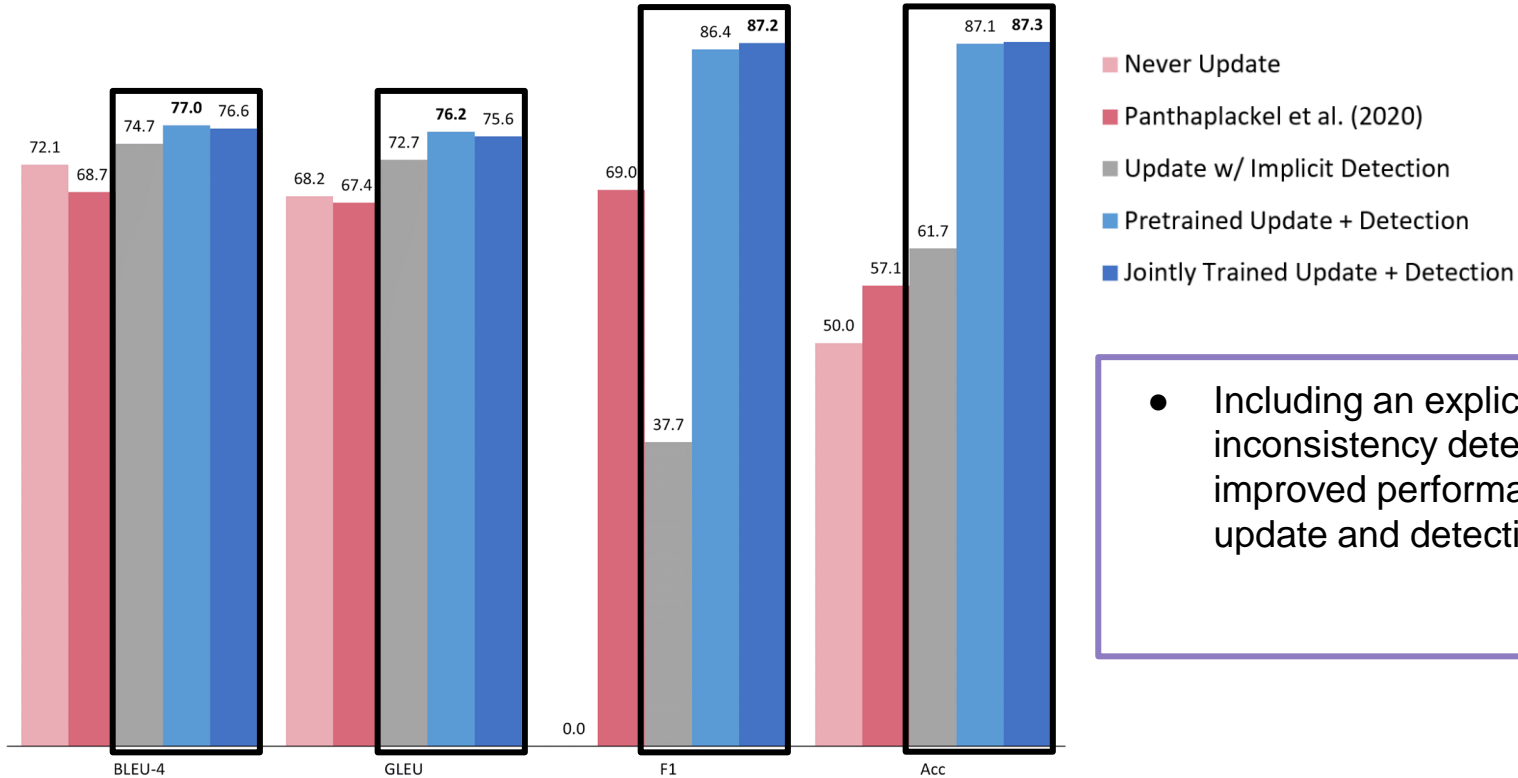
# Extrinsic Evaluation: Results

Cleaned Test Sample



# Extrinsic Evaluation: Results

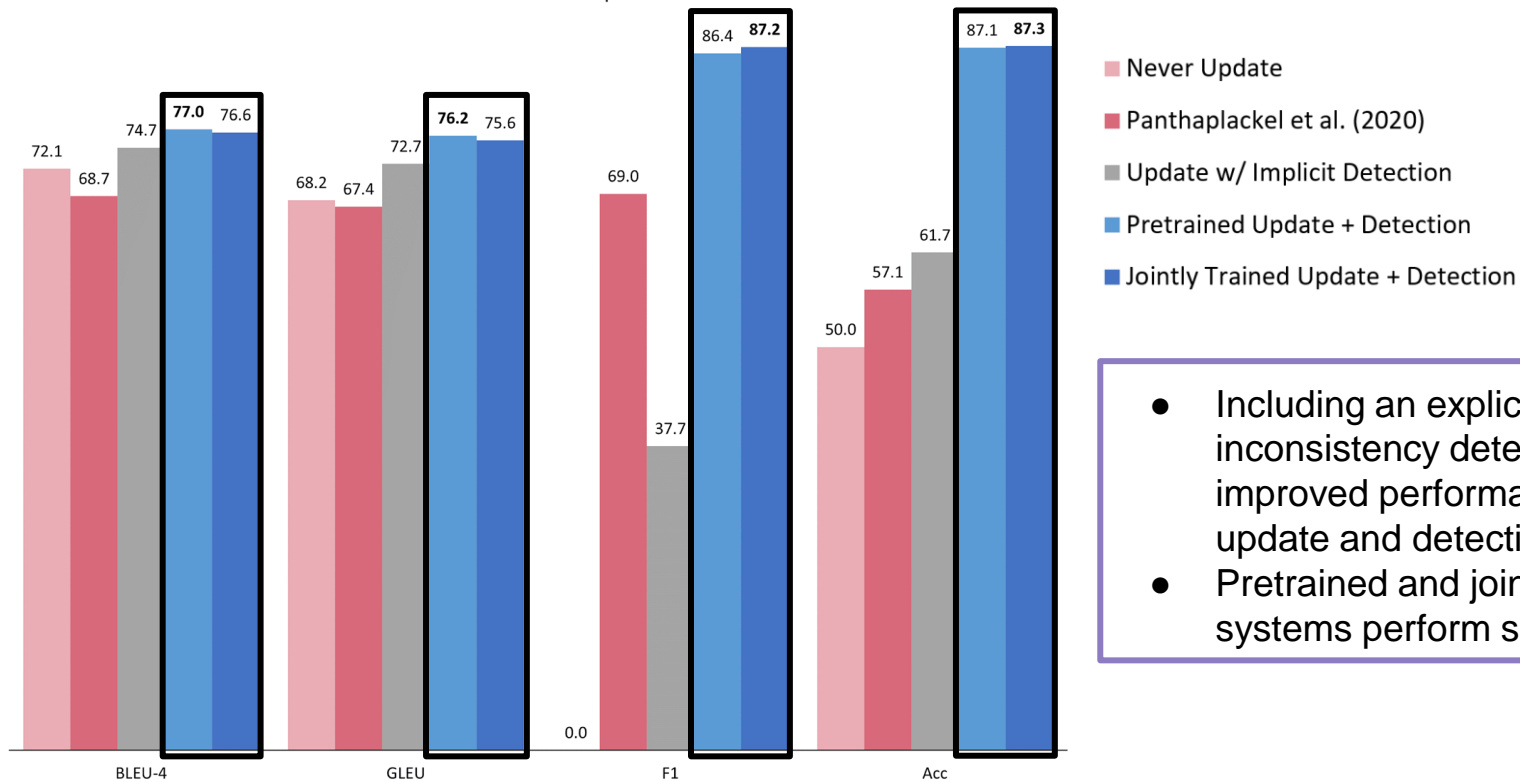
Cleaned Test Sample



- Including an explicit inconsistency detection leads to improved performance across update and detection metrics

# Extrinsic Evaluation: Results

Cleaned Test Sample



- Including an explicit inconsistency detection leads to improved performance across update and detection metrics
- Pretrained and jointly trained systems perform similarly

# Outline

- Task
- Architecture
- Data
- Intrinsic Evaluation
- Extrinsic Evaluation
- **Summary**

# Summary

- We formulated a deep learning approach for just-in-time inconsistency detection between comments and code by learning to relate comments and code changes.
- We show that our approach can outperform multiple baselines as well as post hoc models.
- We also demonstrate that our approach can be used to build a comprehensive comment maintenance system which detects and updates inconsistent comments.

**Code and data available:** <https://github.com/panthap2/deep-jit-inconsistency-detection>

**Contact:** Sheena Panthaplackel <[spantha@cs.utexas.edu](mailto:spantha@cs.utexas.edu)>