

# Producing Explainable AI with Large Language Models

Matthew Renze

AI Ethics

## Question:

Can we use Large Language Models (LLMs) to explain AI predictions in natural language?

# Background

# Background

Explainable AI (XAI)

# Background

Explainable AI (XAI)

Natural language explanation (NLE)

# Background

Explainable AI (XAI)

Natural language explanation (NLE)

Large Language Model (LLM)

# Background

Explainable AI (XAI)

Natural language explanation (NLE)

Large Language Model (LLM)

Generative Pretrained Transformer (GPT)

# Background

Explainable AI (XAI)

Natural language explanation (NLE)

Large Language Model (LLM)

Generative Pretrained Transformer (GPT)

COMPAS Synthetic Database



# Background

Explainable AI (XAI)

Natural language explanation (NLE)

Large Language Model (LLM)

Generative Pretrained Transformer (GPT)

COMPAS Synthetic Database

SHapley Additive exPlanations (SHAP)

# Scaled SHAP Values

c\_charge\_degree\_importance: 0.000199

criminal\_attitude\_importance: 0.000004

family\_criminality\_importance: 0.000002

...

criminal\_associates\_importance: -0.062260

substance\_abuse\_importance: -0.099101

noncompliance\_importance: -0.109318

# Methods

# Experiment

# Experiment

Sample 100 case records

# Experiment

Sample 100 case records

Generate explanations

# Experiment

Sample 100 case records

Generate explanations

Correct explanations

# Experiment

Sample 100 case records

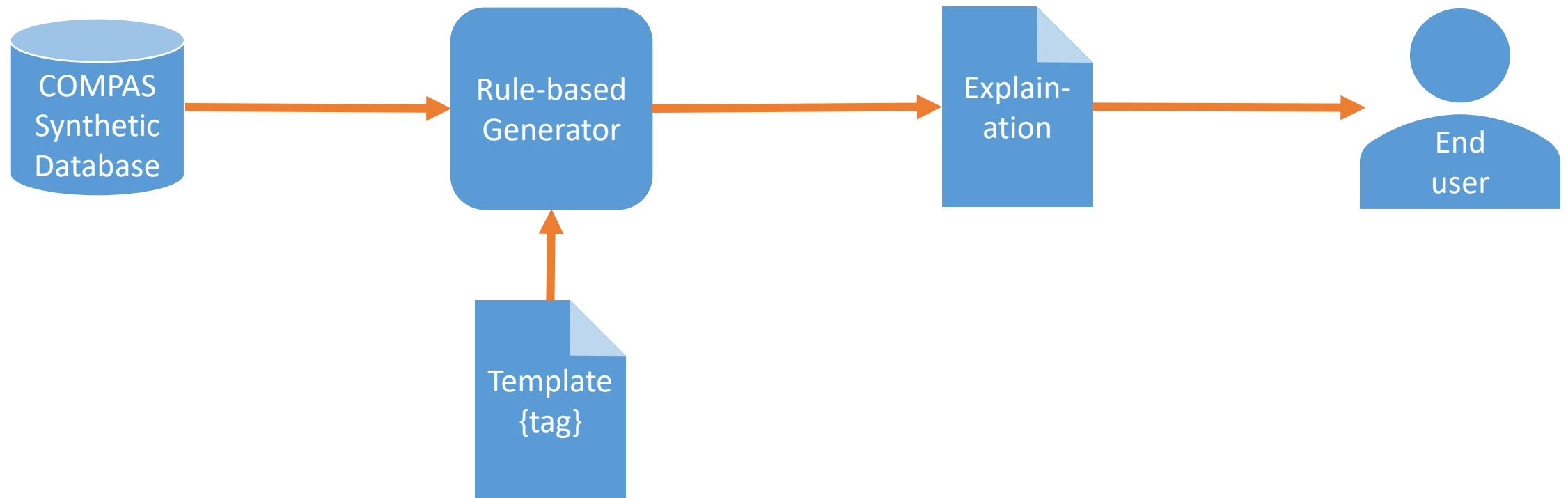
Generate explanations

Correct explanations

Evaluate performance



# Rule-based NLEs



{first-name},

Thank you for sharing your case record for our review.

Regarding your case record #{case-id}, you had a COMPAS risk score of {decile\_score} on a scale of 1-10. This means you were scored as a {decile-score-level} risk for reoffending.

You received this score primarily because you had a {feature-level} risk score for {feature-name} ({feature-score}) {and | ,} {...}. {This was | These were} the {high-importance-feature-count} most important factors in receiving a {decile-score-level} overall risk score.

You also had a {feature-level} risk score for {feature-name} ({feature-score}) {and | ,} {...}. However, {this factor | these factors} only contributed medium importance to the COMPAS algorithm's scoring of your overall risk.

The remaining risk factors including {feature-name} ({feature-score}) {and | ,} {...}, etc. were of low importance to the COMPAS algorithm's score.

Demographic data including age, sex, race, etc., were not included in the algorithm's decision because they are protected by law.

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

{first-name},

Thank you for sharing your case record for our review.

Regarding your case record #{case-id}, you had a COMPAS risk score of {decile\_score} on a scale of 1-10. This means you were scored as a {decile-score-level} risk for reoffending.

You received this score primarily because you had a {feature-level} risk score for {feature-name} ({feature-score}) {and | ,} {...}. {This was | These were} the {high-importance-feature-count} most important factors in receiving a {decile-score-level} overall risk score.

You also had a {feature-level} risk score for {feature-name} ({feature-score}) {and | ,} {...}. However, {this factor | these factors} only contributed medium importance to the COMPAS algorithm's scoring of your overall risk.

The remaining risk factors including {feature-name} ({feature-score}) {and | ,} {...}, etc. were of low importance to the COMPAS algorithm's score.

Demographic data including age, sex, race, etc., were not included in the algorithm's decision because they are protected by law.

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

Richard,

Thank you for sharing your case record for our review.

Regarding your case record #142, you had a COMPAS risk score of 8 on a scale of 1-10. This means you were scored as a high risk for reoffending.

You received this score primarily because you had a high risk score for noncompliance (6) and a high risk score for substance abuse (8). These were the two most important factors in receiving a high overall risk score.

You also had a high risk score for criminal associates (7) and a high risk score for number of prior offenses (9). However, these factors only contributed medium importance to the COMPAS algorithm's scoring of your overall risk.

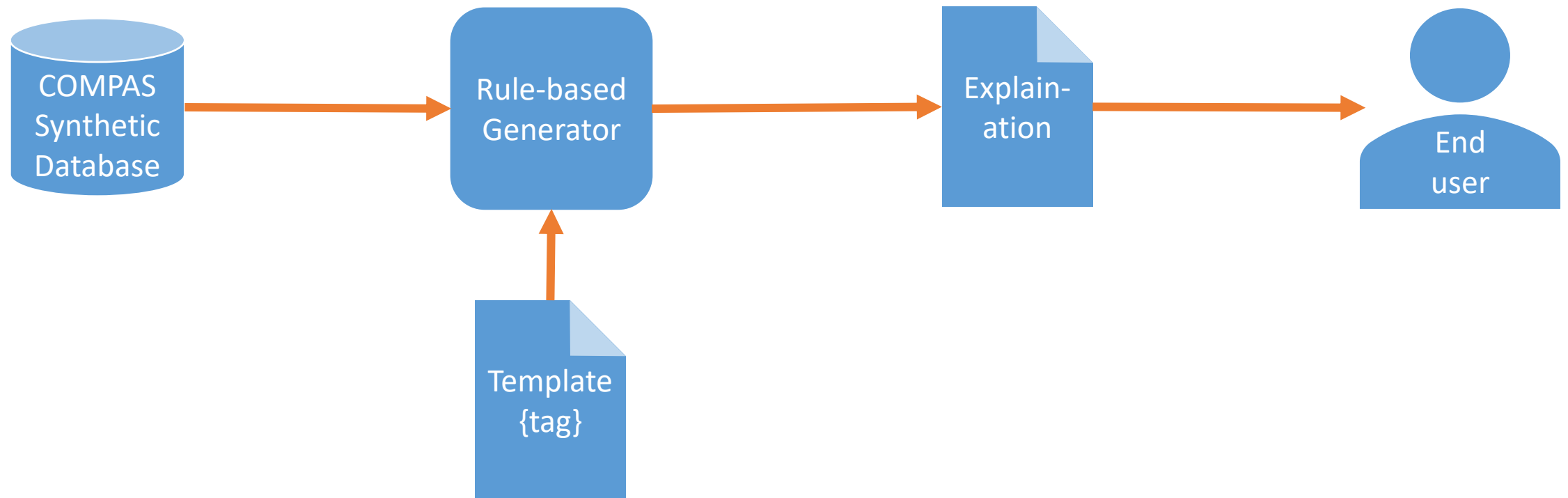
The remaining risk factors including financial problems (3), vocational risk factors (4), and social environment risk (5), etc. were of low importance to the COMPAS algorithm's score.

Demographic data including age, sex, race, etc., were not included in the algorithm's decision because they are protected by law.

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

# Rule-based NLEs



# LLM-based NLEs

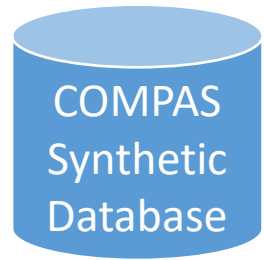


GPT  
Explainer

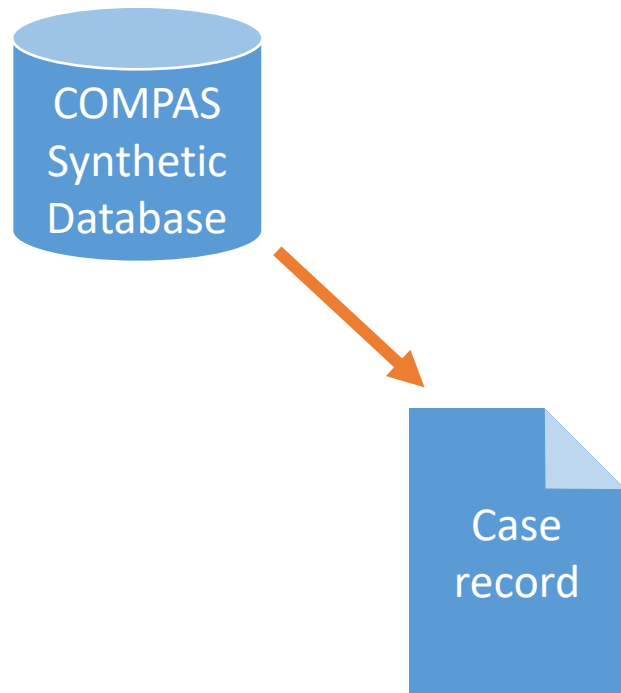
GPT  
Verifier

GPT  
Corrector

# LLM-based NLEs



# LLM-based NLEs





#### # Case Info

id: 142

name: richard veach

decile\_score: 8

decile\_score\_level: high

#### # Features

priors\_count: 9

c\_charge\_degree: F

juv\_fel\_count: 0

family\_criminality: 5

criminal\_attitude: 8

criminal\_associates: 7

financial\_problems: 3

substance\_abuse: 8

noncompliance: 6

social\_environment: 5

vocational: 4

#### # Case Info

id: 142  
name: richard veach  
decile\_score: 8  
decile\_score\_level: high

#### # Features

priors\_count: 9  
c\_charge\_degree: F  
juv\_fel\_count: 0  
family\_criminality: 5  
criminal\_attitude: 8  
criminal\_associates: 7  
financial\_problems: 3  
substance\_abuse: 8  
noncompliance: 6  
social\_environment: 5  
vocational: 4

#### # Feature Levels

priors\_count\_level: high  
c\_charge\_degree\_level: high  
juv\_fel\_count\_level: low  
family\_criminality\_level: medium  
criminal\_attitude\_level: high  
criminal\_associates\_level: high  
financial\_problems\_level: low  
substance\_abuse\_level: high  
noncompliance\_level: high  
social\_environment\_level: medium  
vocational\_level: medium

### # Case Info

id: 142  
name: richard veach  
decile\_score: 8  
decile\_score\_level: high

### # Features

priors\_count: 9  
c\_charge\_degree: F  
juv\_fel\_count: 0  
family\_criminality: 5  
criminal\_attitude: 8  
criminal\_associates: 7  
financial\_problems: 3  
substance\_abuse: 8  
noncompliance: 6  
social\_environment: 5  
vocational: 4

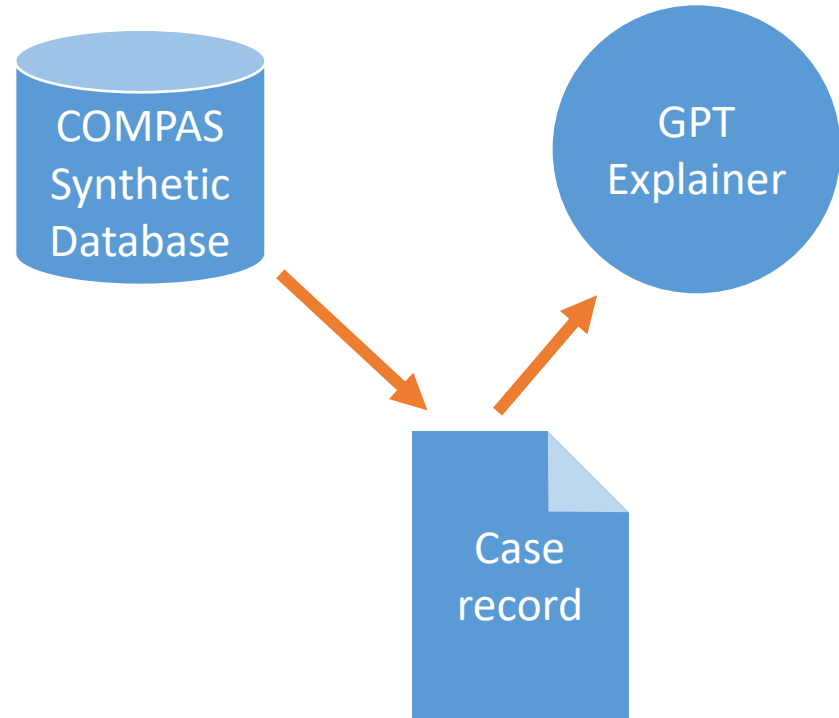
### # Feature Levels

priors\_count\_level: high  
c\_charge\_degree\_level: high  
juv\_fel\_count\_level: low  
family\_criminality\_level: medium  
criminal\_attitude\_level: high  
criminal\_associates\_level: high  
financial\_problems\_level: low  
substance\_abuse\_level: high  
noncompliance\_level: high  
social\_environment\_level: medium  
vocational\_level: medium

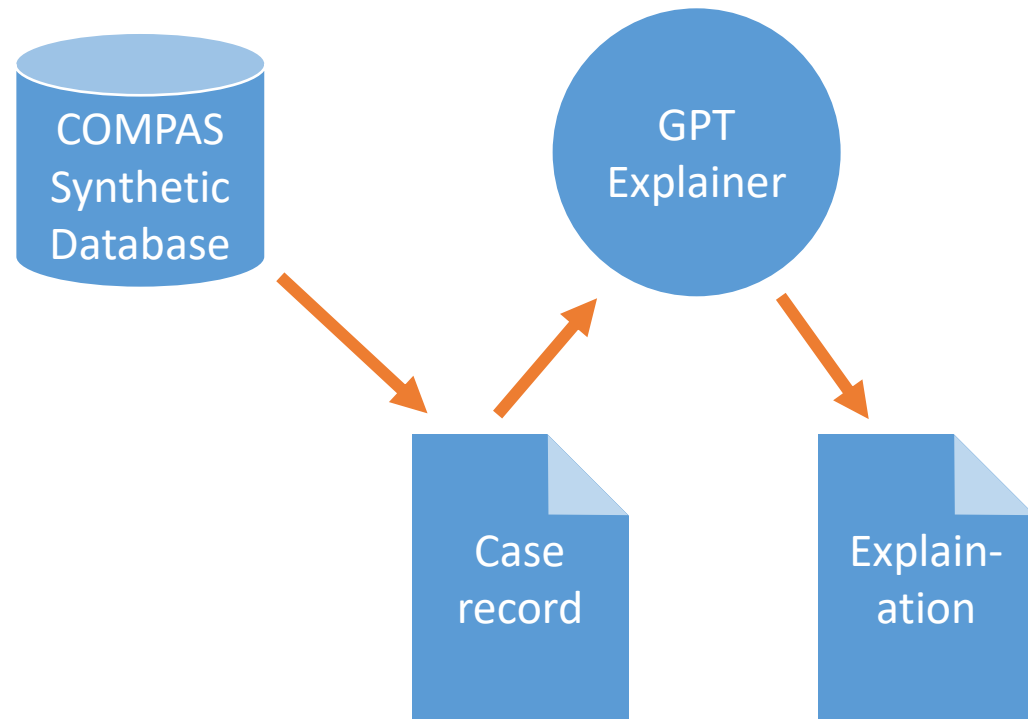
### # Feature Importance Levels

noncompliance\_importance\_level: high  
substance\_abuse\_importance\_level: high  
criminal\_associates\_importance\_level: medium  
priors\_count\_importance\_level: medium  
financial\_problems\_importance\_level: low  
vocational\_importance\_level: low  
social\_environment\_importance\_level: low  
c\_charge\_degree\_importance\_level: low  
juv\_fel\_count\_importance\_level: low  
family\_criminality\_importance\_level: low  
criminal\_attitude\_importance\_level: low

# LLM-based NLEs



# LLM-based NLEs



Richard,

Thank you for sharing your case record for our review.

Regarding your case record #142, you had a COMPAS risk score of 8 on a scale of 1-10. This means you were scored as a high risk for reoffending.

You received this score primarily because you had a high risk score for criminal attitude (8), a high risk score for noncompliance (6), and a high risk score for substance abuse (8). These were the three most important factors in receiving a high overall risk score.

You also had a high risk score for criminal associates (7) and a medium risk score for number of prior offenses (9). However, these factors only contributed medium importance to the COMPAS algorithm's scoring of your overall risk.

The remaining risk factors including financial problems (3), social environment risk (5), and vocational risk factors (4), etc. were of low importance to the COMPAS algorithm's score.

Demographic data including age, sex, race, etc., were not included in the algorithm's decision because they are protected by law.

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

Richard,

Thank you for sharing your case record for our review.

Regarding your case record #142, you had a COMPAS risk score of 8 on a scale of 1-10. This means you were scored as a high risk for reoffending.

You received this score primarily because you had a high risk score for **criminal attitude (8)**, a high risk score for noncompliance (6), and a high risk score for substance abuse (8). These were the three most important factors in receiving a high overall risk score.

You also had a high risk score for criminal associates (7) and a **medium** risk score for number of prior offenses (9). However, these factors only contributed medium importance to the COMPAS algorithm's scoring of your overall risk.

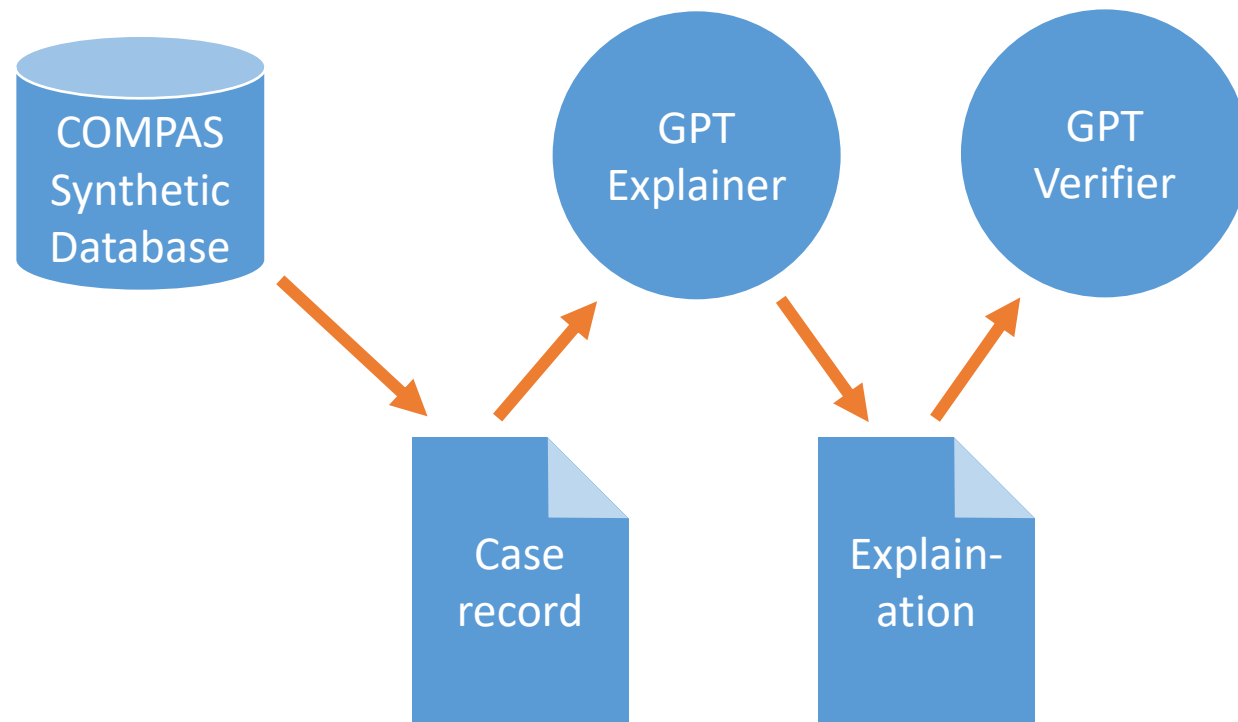
The remaining risk factors including financial problems (3), social environment risk (5), and vocational risk factors (4), etc. were of low importance to the COMPAS algorithm's score.

Demographic data including age, sex, race, etc., were not included in the algorithm's decision because they are protected by law.

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

# LLM-based NLEs





Richard,

Thank you for sharing your case record for our review.

Regarding your case record #142, you had a COMPAS risk score of 8 on a scale of 1-10. This means you were scored as a high risk for reoffending.

You received this score primarily because you had a high risk score for criminal attitude (8), a high risk score for noncompliance (6), and a high risk score for substance abuse (8). These were the three most important factors in receiving a high overall risk score.

You also had a high risk score for criminal associates (7) and a medium risk score for number of prior offenses (9). However, these factors only contributed medium importance to the COMPAS algorithm's scoring of your overall risk.

The remaining risk factors including financial problems (3), social environment risk (5), and vocational risk factors (4), etc. were of low importance to the COMPAS algorithm's score.

Demographic data including age, sex, race, etc., were not included in the algorithm's decision because they are protected by law.

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

Richard,

Thank you for sharing your case record for our review.

Regarding your case record #142, you had a COMPAS risk score of 8 on a scale of 1-10. This means you were scored as a high risk for reoffending.

You received this score primarily because you had a high risk score for criminal attitude (8), a high risk score for noncompliance (6), and a high risk score for substance abuse (8). These were the three most important factors in receiving a high overall risk score.

You also had a high risk score for criminal associates (7) and a medium risk score for number of prior offenses (9). However, these factors only contributed medium importance to the COMPAS algorithm's scoring of your overall risk.

The remaining risk factors including financial problems (3), social environment risk (5), and vocational risk factors (4), etc. were of low importance to the COMPAS algorithm's score.

Demographic data including age, sex, race, etc., were not included in the algorithm's decision because they are protected by law.

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

### # Case Info

id: 142  
name: richard  
decile\_score: 8  
decile\_score\_level: high

### # Features

priors\_count: 9  
c\_charge\_degree:  
juv\_fel\_count:  
family\_criminality:  
criminal\_attitude:  
criminal\_associates: 7  
financial\_problems: 3  
substance\_abuse: 8  
noncompliance: 6  
social\_environment: 5  
vocational: 4

### # Feature Levels

priors\_count\_level: high  
c\_charge\_degree\_level:  
juv\_fel\_count\_level:  
family\_criminality\_level:  
criminal\_attitude\_level:  
criminal\_associates\_level: high  
financial\_problems\_level: low  
substance\_abuse\_level: high  
noncompliance\_level: high  
social\_environment\_level: medium  
vocational\_level: medium

### # Feature Importance

noncompliance\_importance\_level: high  
substance\_abuse\_importance\_level: high  
criminal\_associates\_importance\_level: medium  
priors\_count\_importance\_level: medium  
financial\_problems\_importance\_level: low  
social\_environment\_importance\_level: low  
c\_charge\_degree\_importance\_level: low  
juv\_fel\_count\_importance\_level: low  
family\_criminality\_importance\_level: low  
criminal\_attitude\_importance\_level: low  
vocational\_importance\_level: low

### # Case Info

id: 142  
name: richard  
decile\_score: 8  
decile\_score\_level: high

### # Features

priors\_count: 9  
c\_charge\_degree:   
juv\_fel\_count:   
family\_criminality:   
criminal\_attitude:   
criminal\_associates: 7  
financial\_problems: 3  
substance\_abuse: 8  
noncompliance: 6  
social\_environment: 5  
vocational: 4

### # Feature Levels

priors\_count\_level: high  
c\_charge\_degree\_level:   
juv\_fel\_count\_level:   
family\_criminality\_level:   
criminal\_attitude\_level:   
criminal\_associates\_level: high  
financial\_problems\_level: low  
substance\_abuse\_level: high  
noncompliance\_level: high  
social\_environment\_level: medium  
vocational\_level: medium

### # Feature Importance

noncompliance\_importance\_level: high  
substance\_abuse\_importance\_level: high  
criminal\_associates\_importance\_level: medium  
priors\_count\_importance\_level: medium  
financial\_problems\_importance\_level: low  
social\_environment\_importance\_level: low  
c\_charge\_degree\_importance\_level: low  
juv\_fel\_count\_importance\_level: low  
family\_criminality\_importance\_level: low  
criminal\_attitude\_importance\_level: low  
vocational\_importance\_level: low

# Errors

priors\_count\_level: medium -> high

vocational\_level: low -> medium

criminal\_attitude\_importance\_level: high -> low

# Errors

priors\_count\_level: **medium** -> high

vocational\_level: **low** -> medium

criminal\_attitude\_importance\_level: **high** -> low

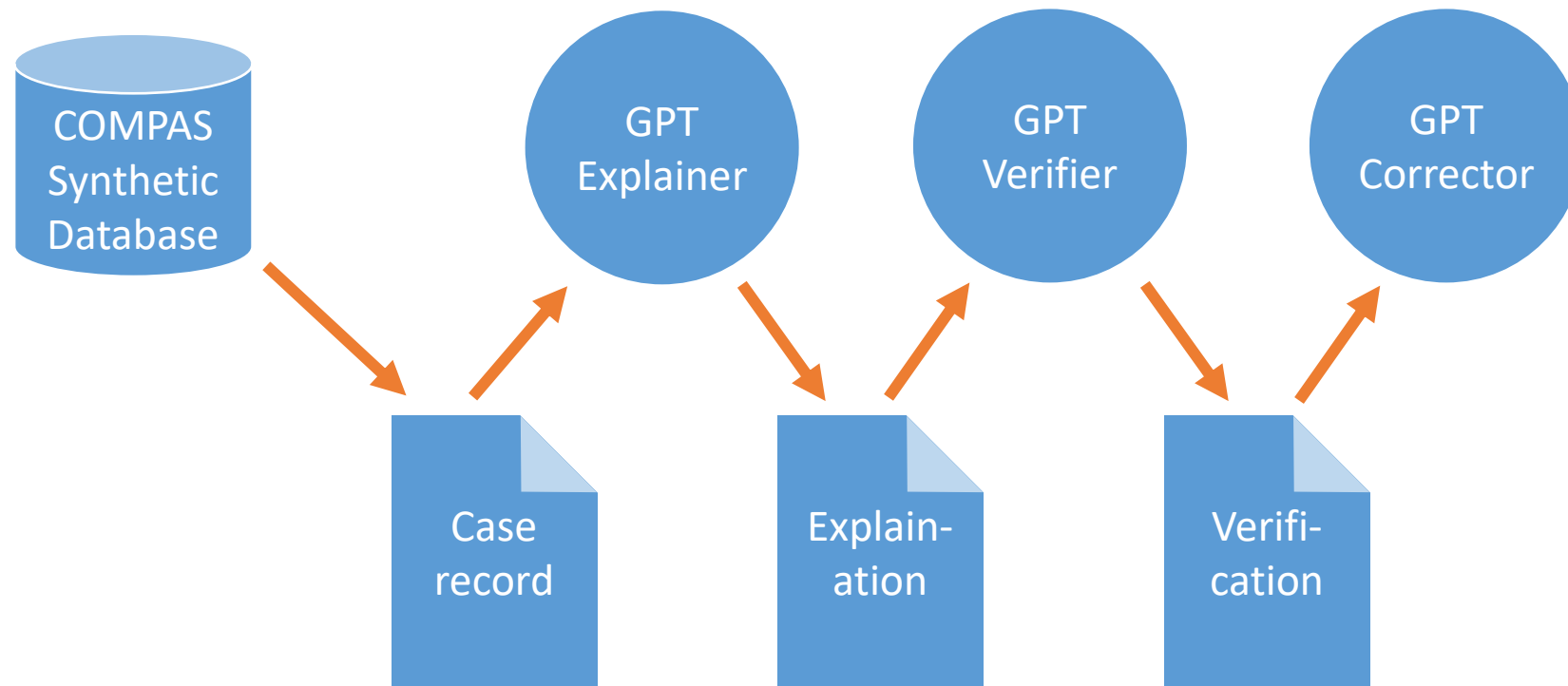
# Errors

priors\_count\_level: **medium** -> **high**

vocational\_level: **low** -> **medium**

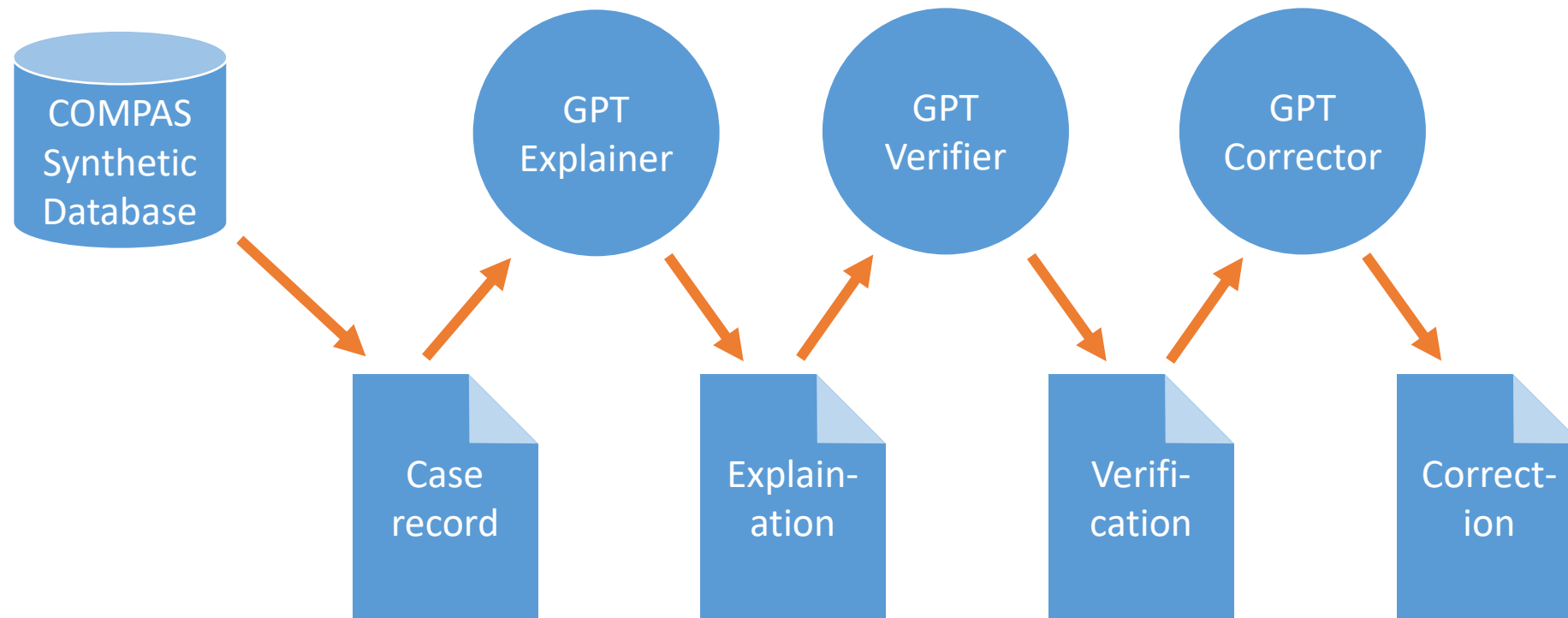
criminal\_attitude\_importance\_level: **high** -> **low**

# LLM-based NLEs





# LLM-based NLEs



Richard,

Thank you for sharing your case record for our review.

Regarding your case record #142, you had a COMPAS risk score of 8 on a scale of 1-10. This means you were scored as a high risk for reoffending.

You received this score primarily because you had a high risk score for **criminal attitude (8)**, a high risk score for noncompliance (6), and a high risk score for substance abuse (8). These were the three most important factors in receiving a high overall risk score.

You also had a high risk score for criminal associates (7) and a **medium** risk score for number of prior offenses (9). However, these factors only contributed medium importance to the COMPAS algorithm's scoring of your overall risk.

The remaining risk factors including financial problems (3), social environment risk (5), and vocational risk factors (4), etc. were of low importance to the COMPAS algorithm's score.

Demographic data including age, sex, race, etc., were not included in the algorithm's decision because they are protected by law.

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

Richard,

Thank you for sharing your case record for our review.

Regarding your case record #142, you had a COMPAS risk score of 8 on a scale of 1-10. This means you were scored as a high risk for reoffending.

You received this score primarily because you had a high risk score for noncompliance (6) and a high risk score for substance abuse (8). These were the two most important factors in receiving a high overall risk score.

You also had a high risk score for criminal associates (7) and a high risk score for the number of prior offenses (9). However, these factors only contributed medium importance to the COMPAS algorithm's scoring of your overall risk.

The remaining risk factors including financial problems (3), social environment risk (5), and vocational risk factors (4), etc. were of low importance to the COMPAS algorithm's score.

Demographic data including age, sex, race, etc., were not included in the algorithm's decision because they are protected by law.

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

Richard,

Thank you for sharing your case record for our review.

...

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

#### # Corrections

- Line 26: Fixed priors\_count\_level: medium -> high
- Line 27: Fixed vocational\_level: low -> medium
- Line 32: Fixed criminal\_attitude\_importance\_level: high -> low

#### # Explanation

- The priors\_count feature had a value of 9, which is high, so the feature level was corrected from medium to high.
- The vocational feature had a value of 4, which is medium, so the feature level was corrected from low to medium.
- The criminal\_attitude feature had a high risk score of 8, but it was incorrectly listed as having high importance. It was corrected to low importance.

Richard,

Thank you for sharing your case record for our review.

...

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

#### # Corrections

- Line 26: Fixed priors\_count\_level: medium -> high
- Line 27: Fixed vocational\_level: low -> medium
- Line 32: Fixed criminal\_attitude\_importance\_level: high -> low

#### # Explanation

- The priors\_count feature had a value of 9, which is high, so the feature level was corrected from medium to high.
- The vocational feature had a value of 4, which is medium, so the feature level was corrected from low to medium.
- The criminal\_attitude feature had a high risk score of 8, but it was incorrectly listed as having high importance. It was corrected to low importance.

Richard,

Thank you for sharing your case record for our review.

...

Please let us know if you have any questions or concerns.

Sincerely,  
Explainable AI

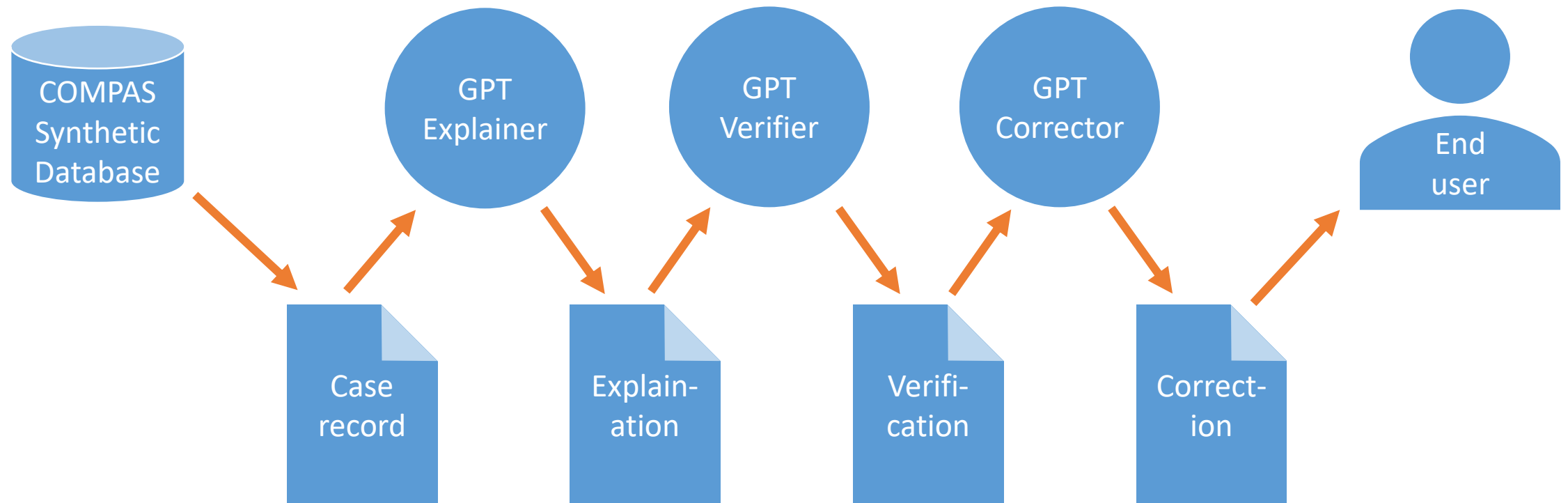
#### # Corrections

- Line 26: Fixed priors\_count\_level: medium -> high
- Line 27: Fixed vocational\_level: low -> medium
- Line 32: Fixed criminal\_attitude\_importance\_level: high -> low

#### # Explanation

- The priors\_count feature had a value of 9, which is high, so the feature level was corrected from medium to high.
- The vocational feature had a value of 4, which is medium, so the feature level was corrected from low to medium.
- The criminal\_attitude feature had a high risk score of 8, but it was incorrectly listed as having high importance. It was corrected to low importance.

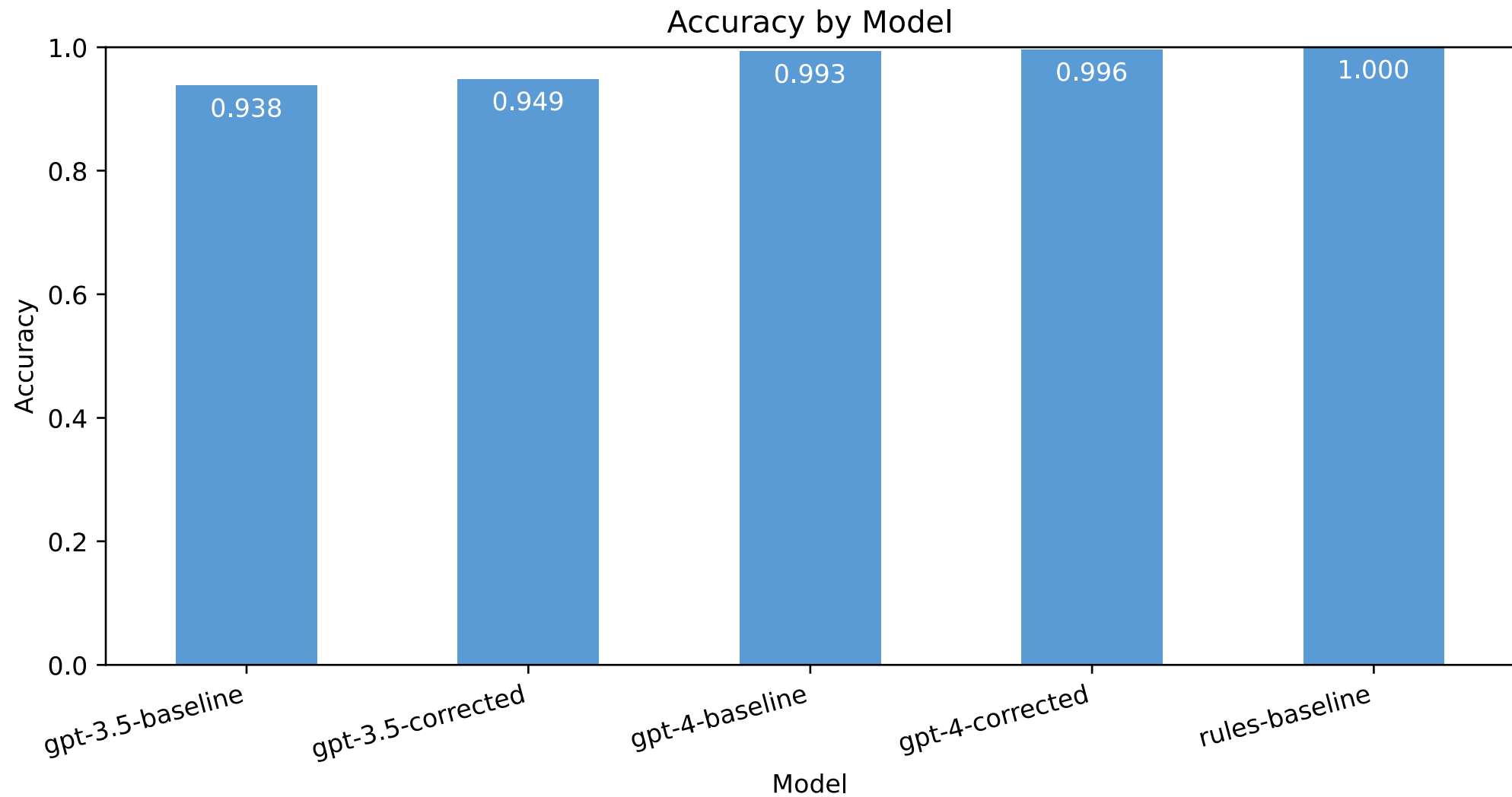
# LLM-based NLEs



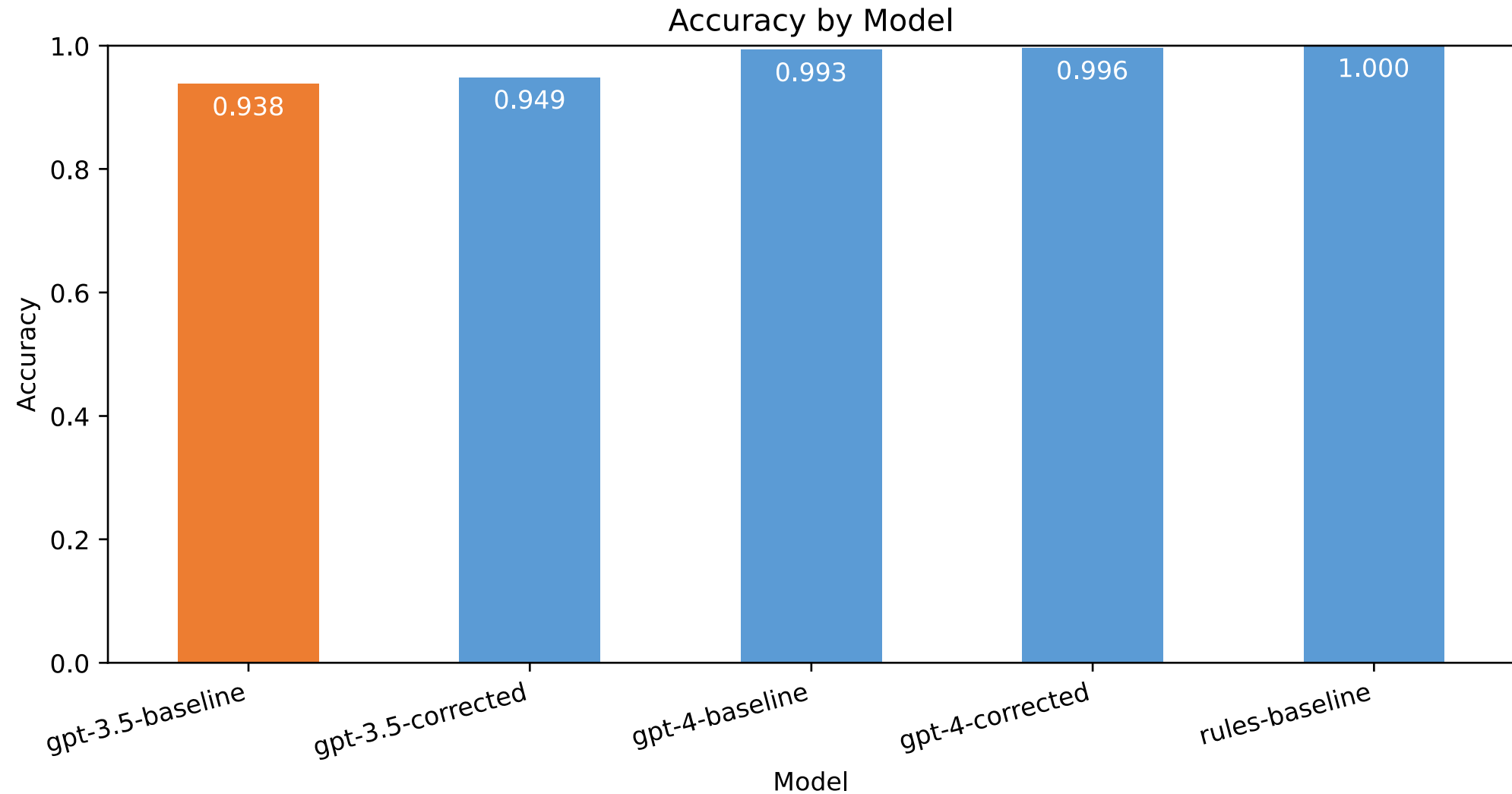
# Results



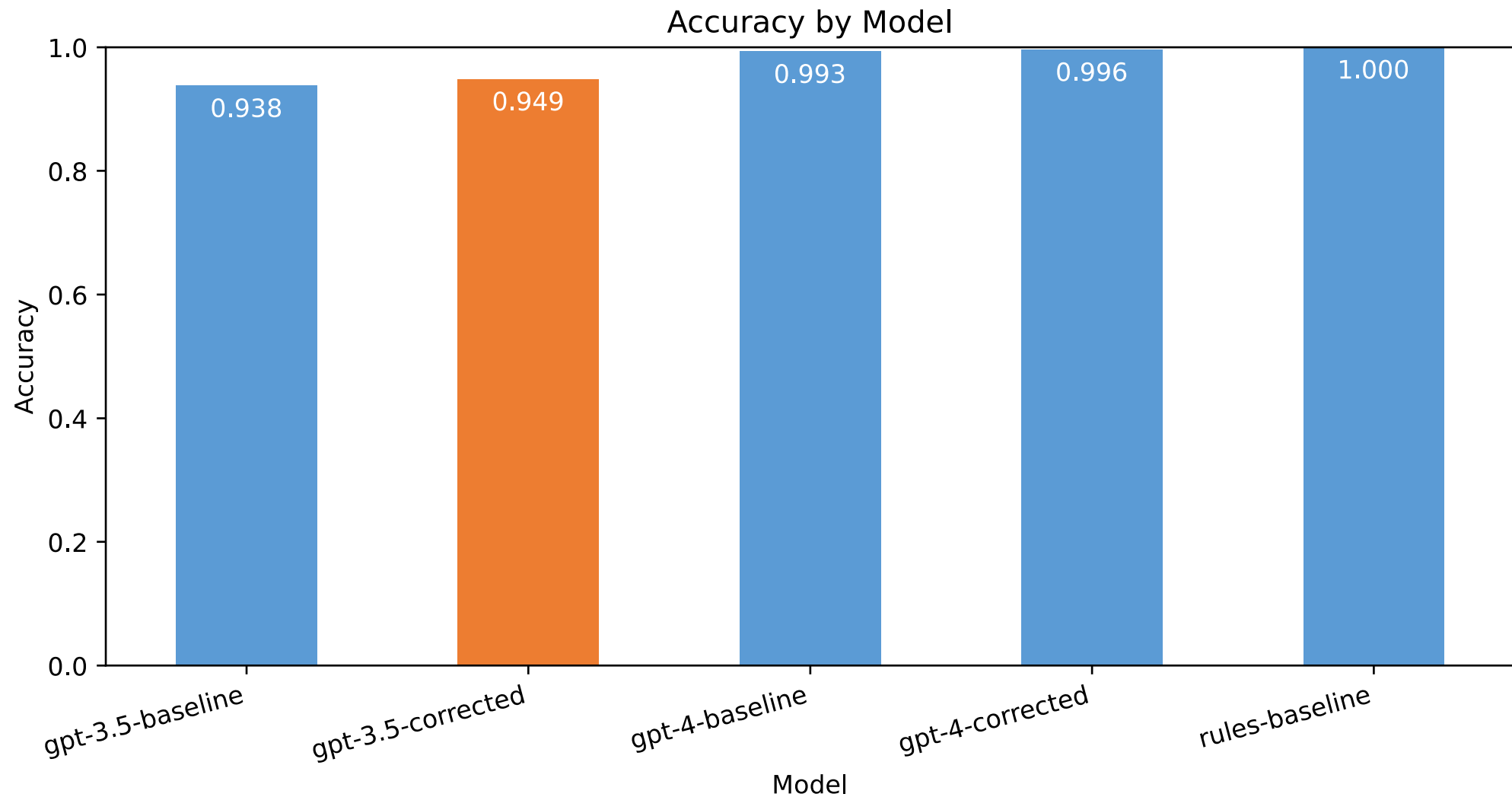
# Factual Accuracy



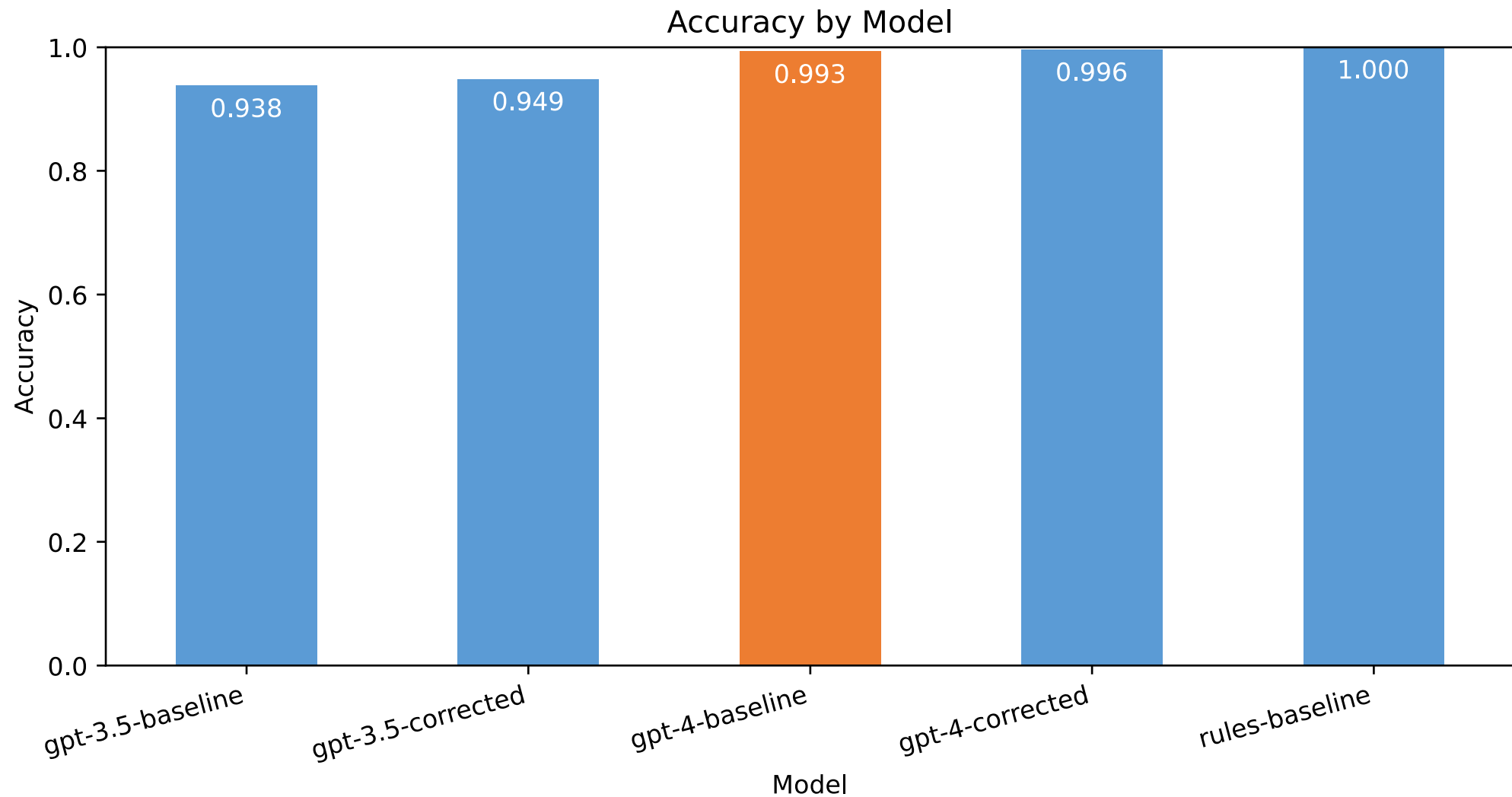
# Factual Accuracy



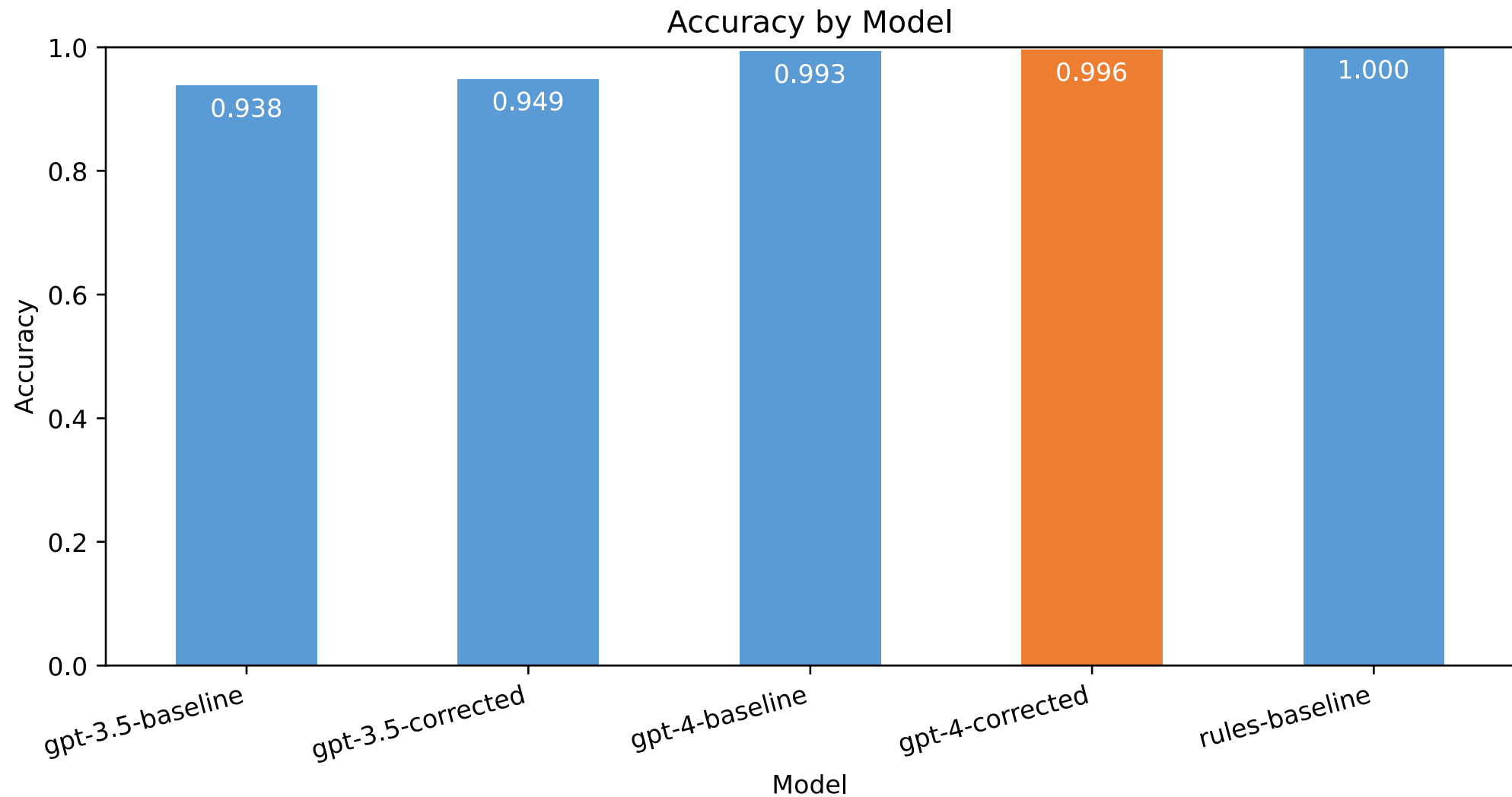
# Factual Accuracy



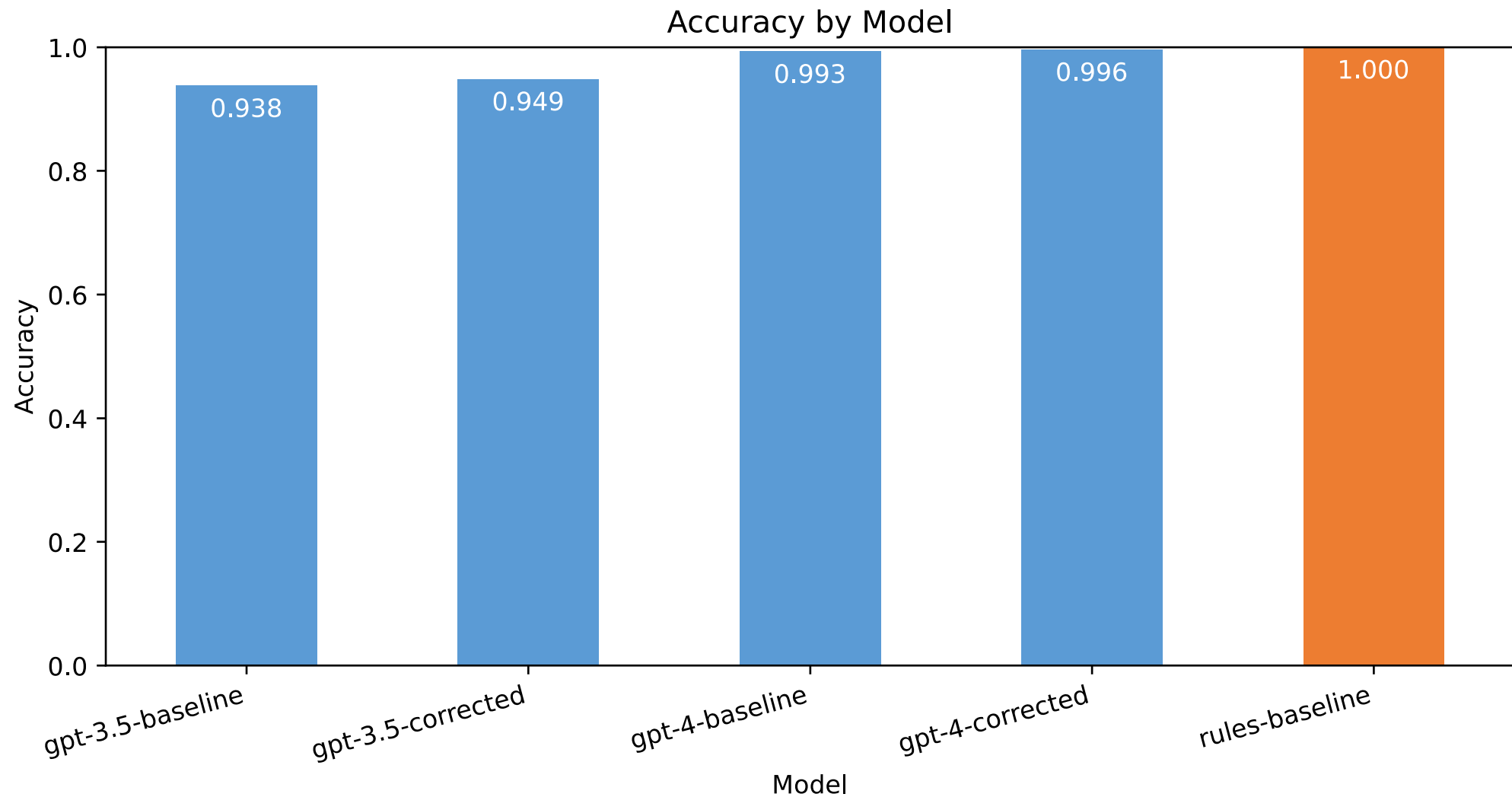
# Factual Accuracy



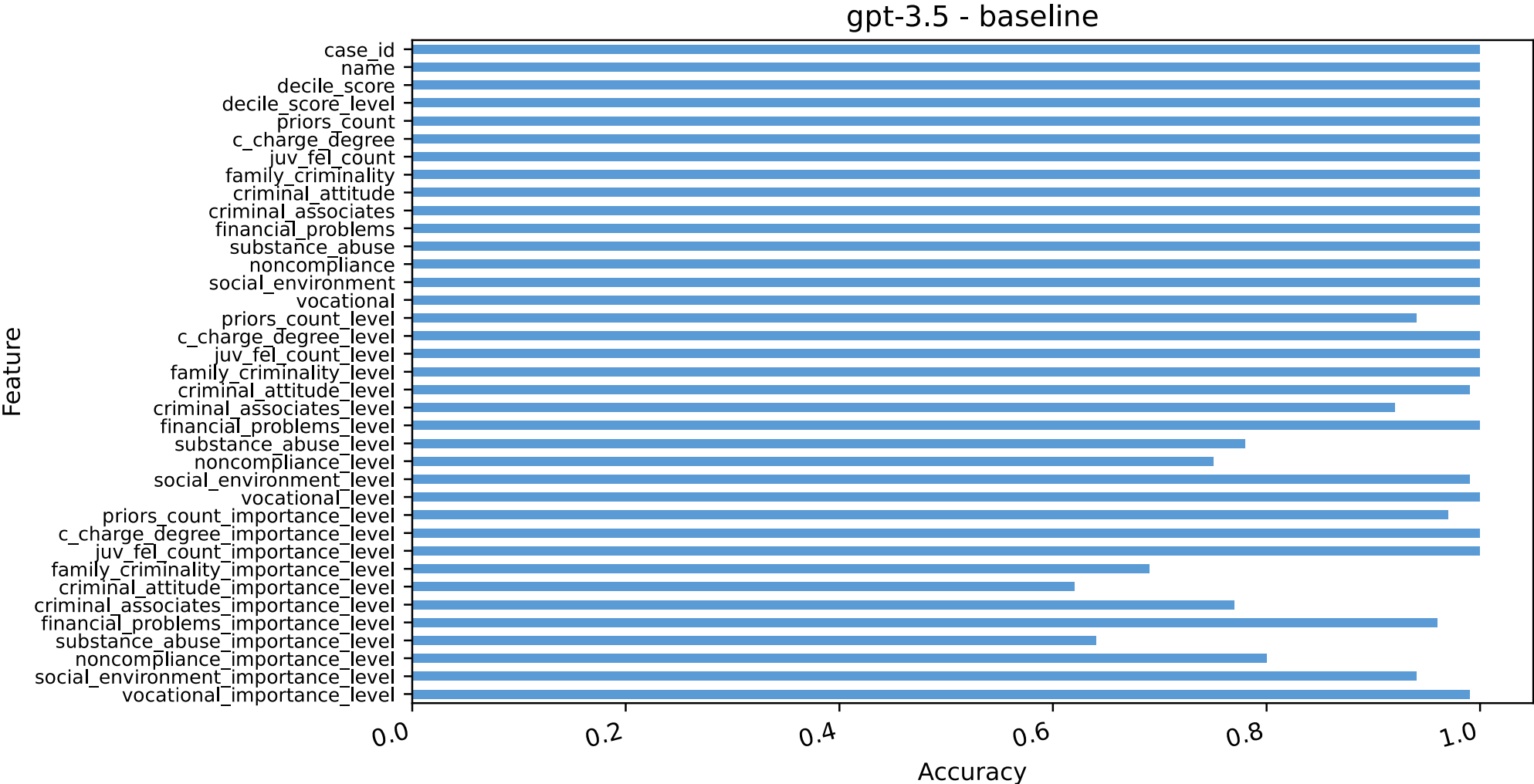
# Factual Accuracy



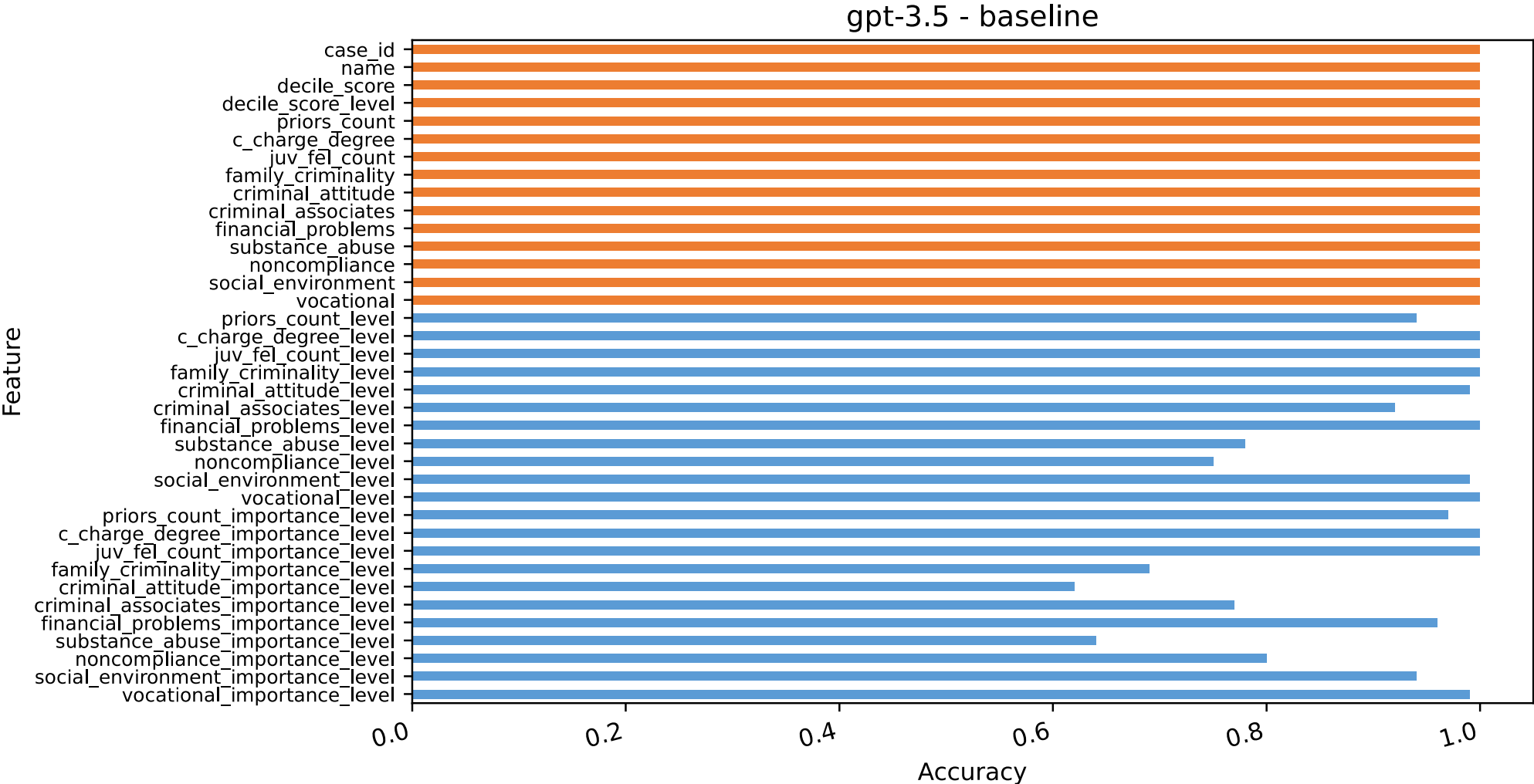
# Factual Accuracy



# Detailed Analysis (Worst Case)

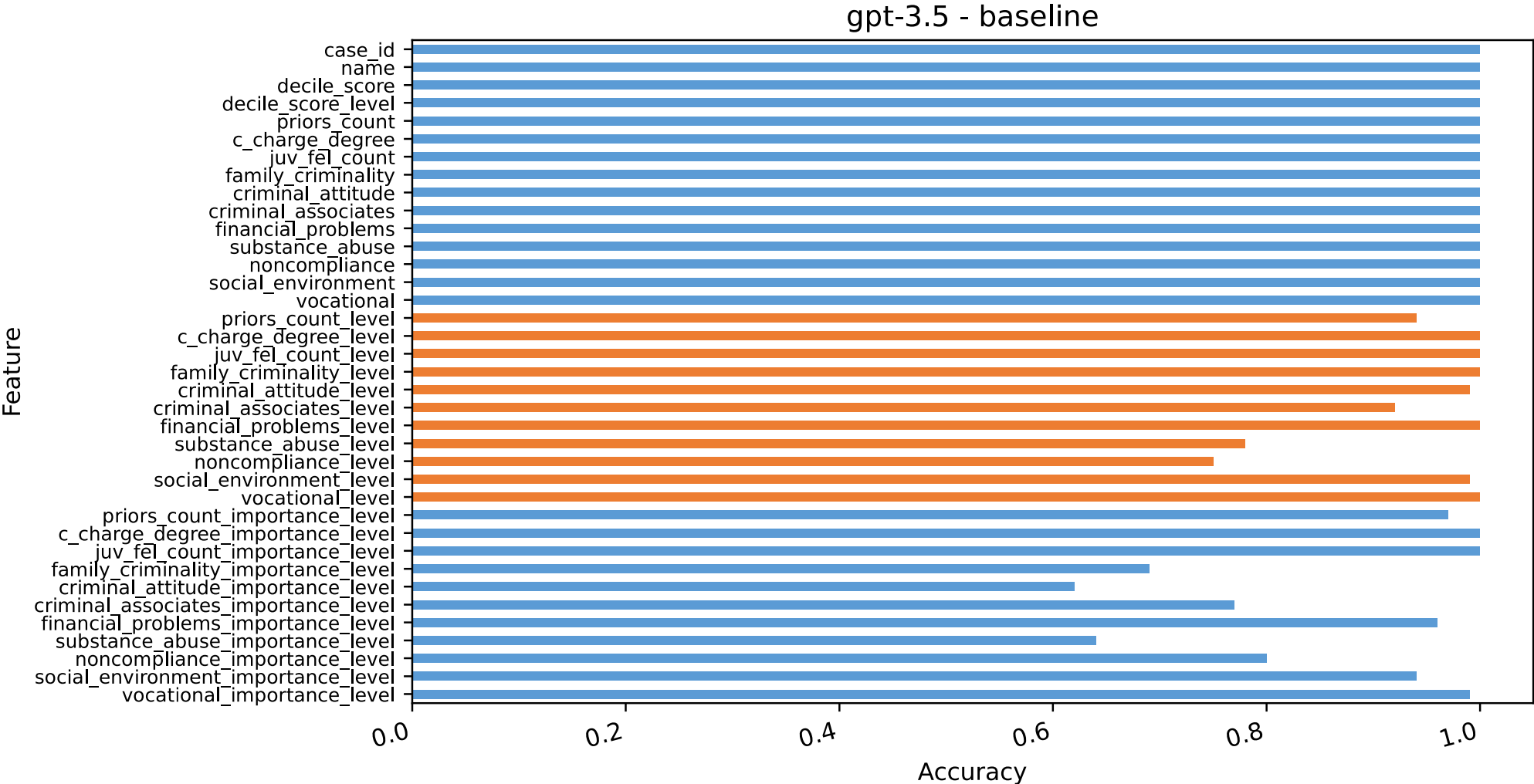


# Detailed Analysis (Worst Case)

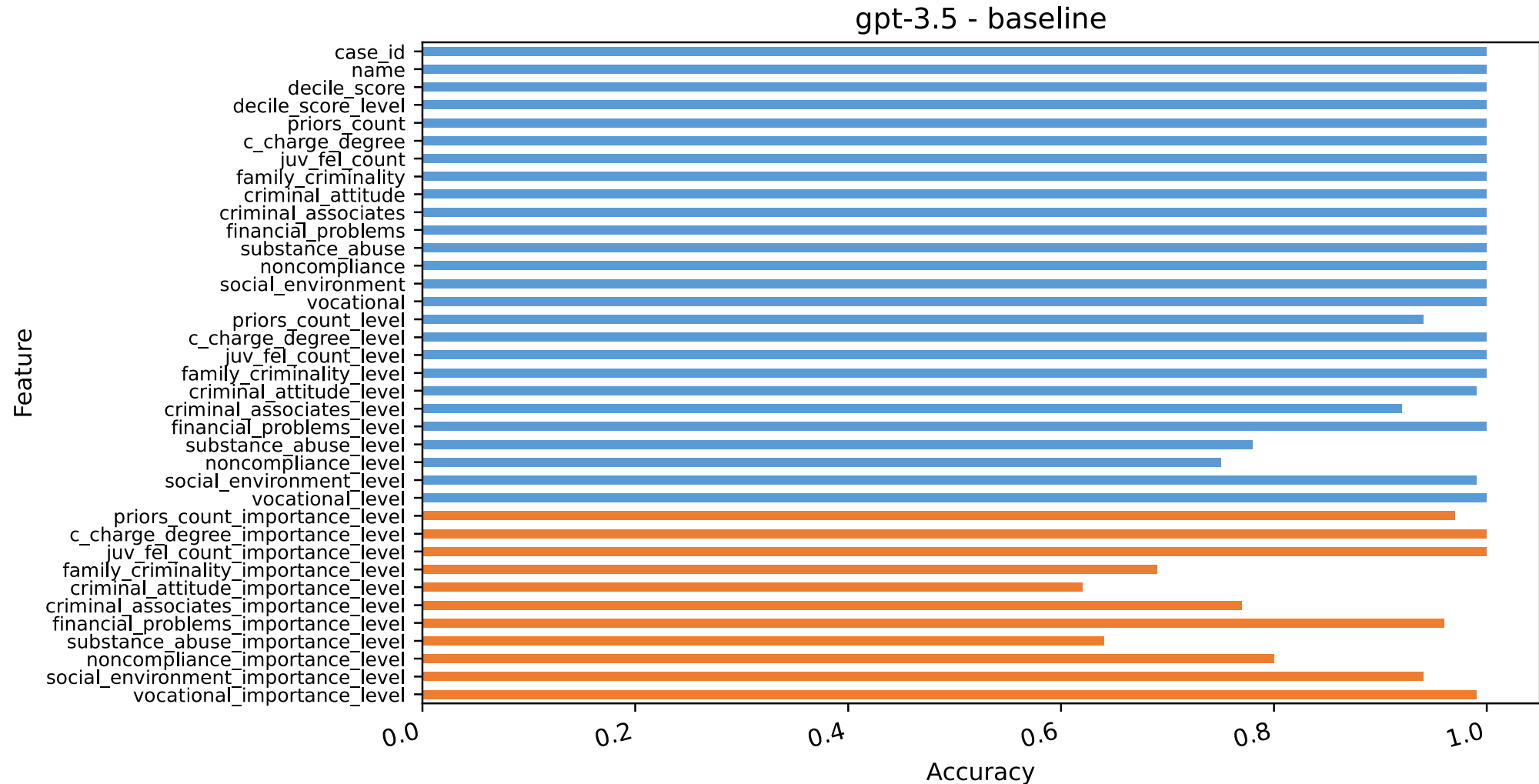




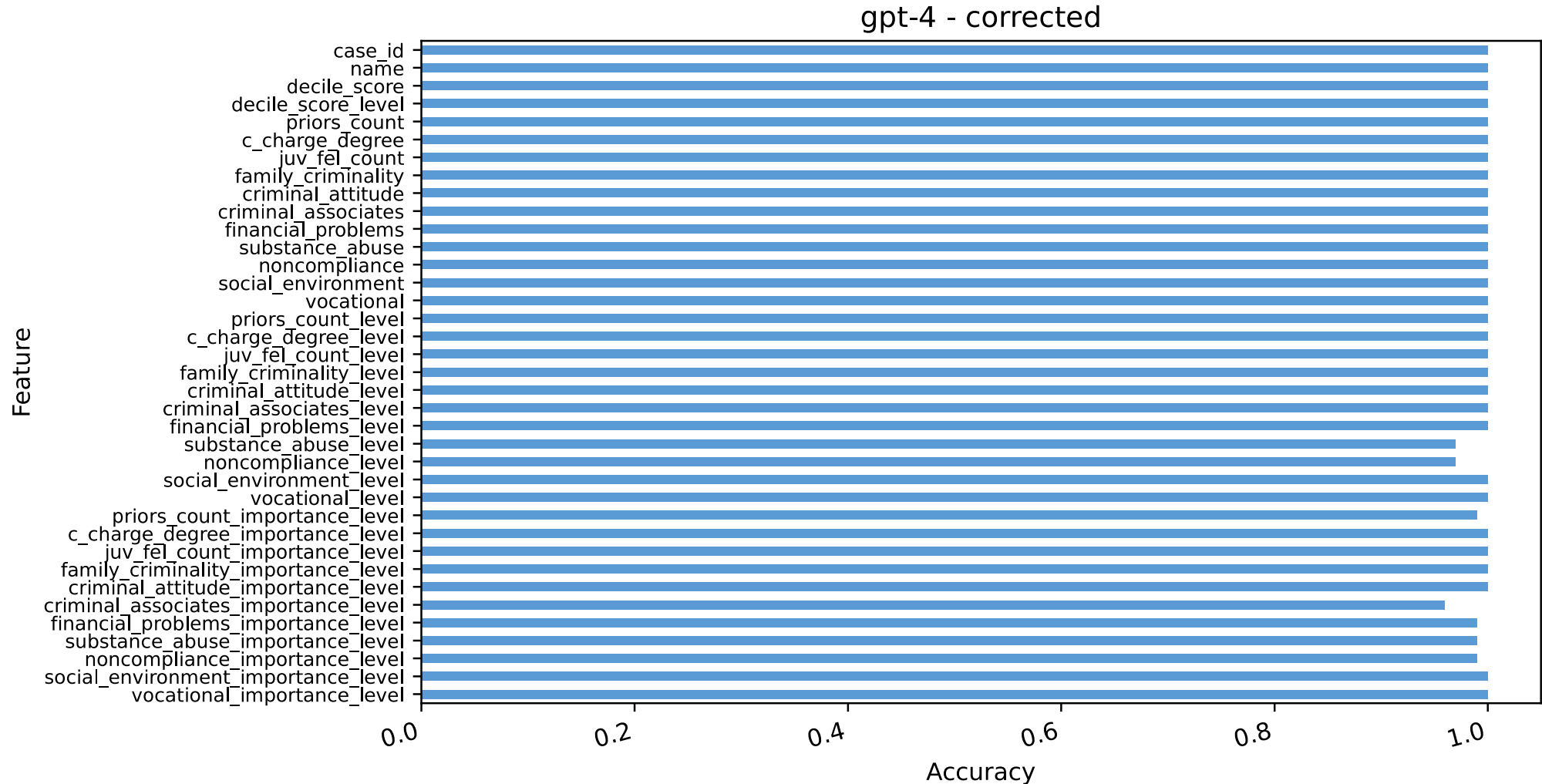
# Detailed Analysis (Worst Case)



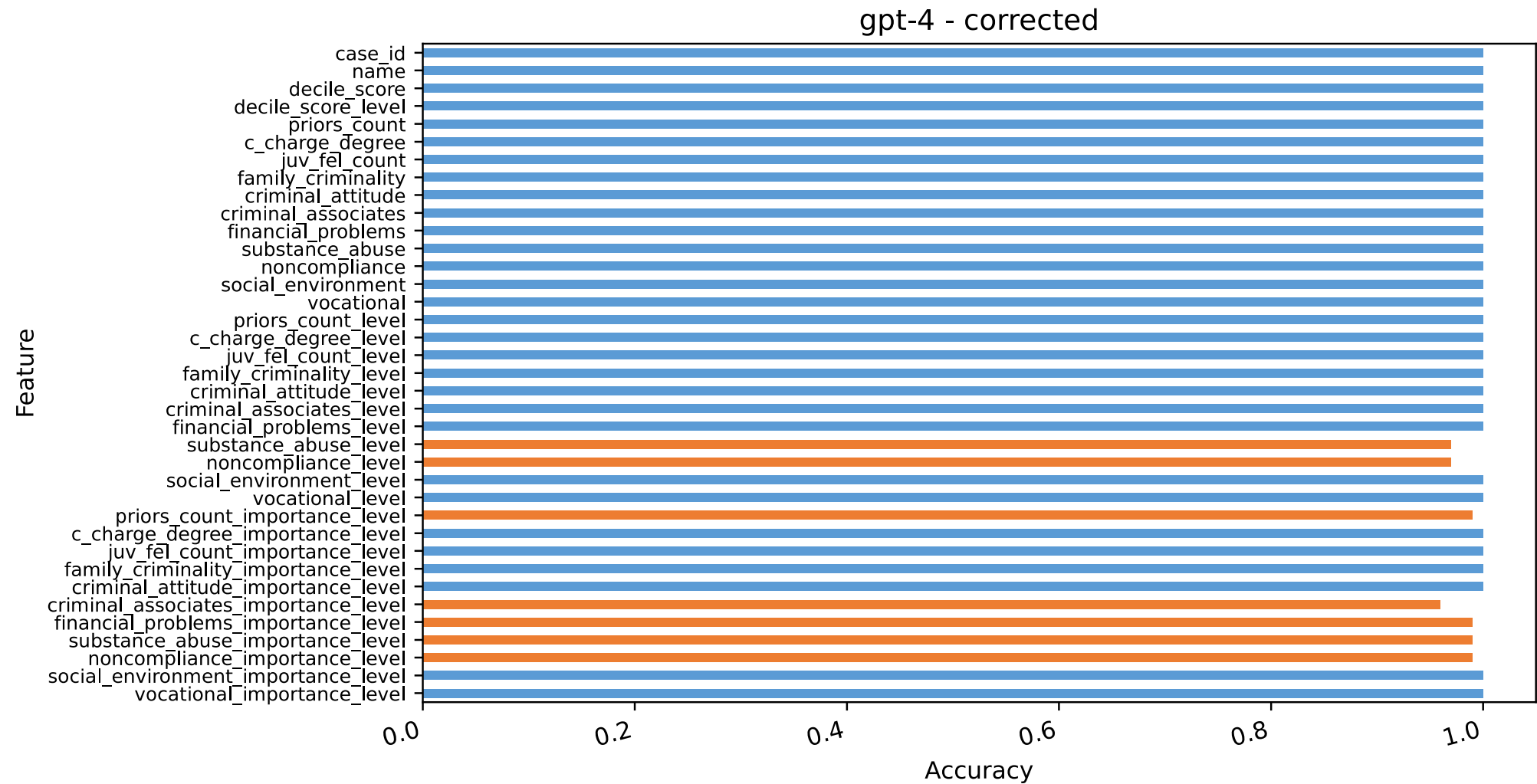
# Detailed Analysis (Worst Case)



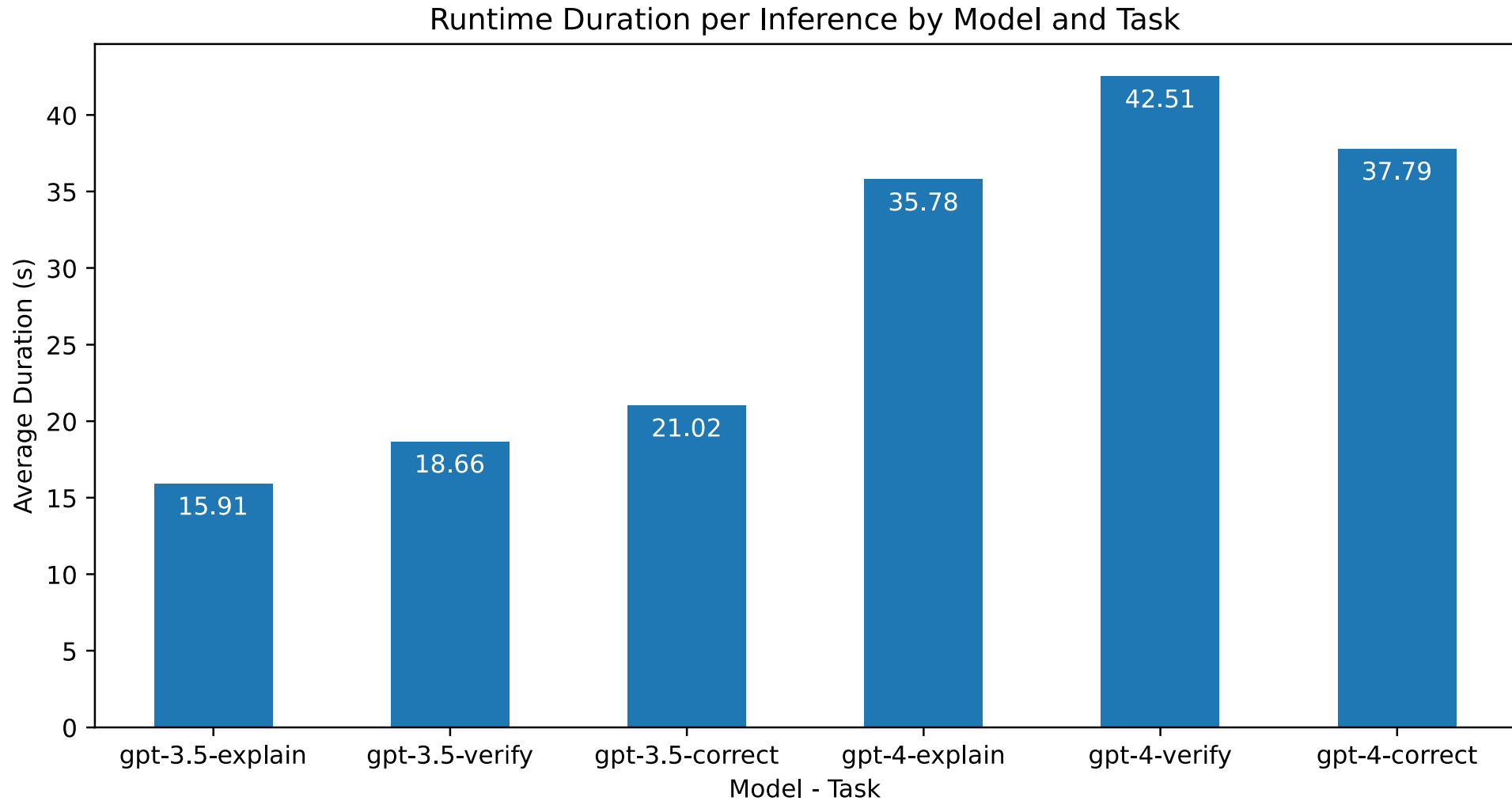
# Detailed Analysis (Best Case)



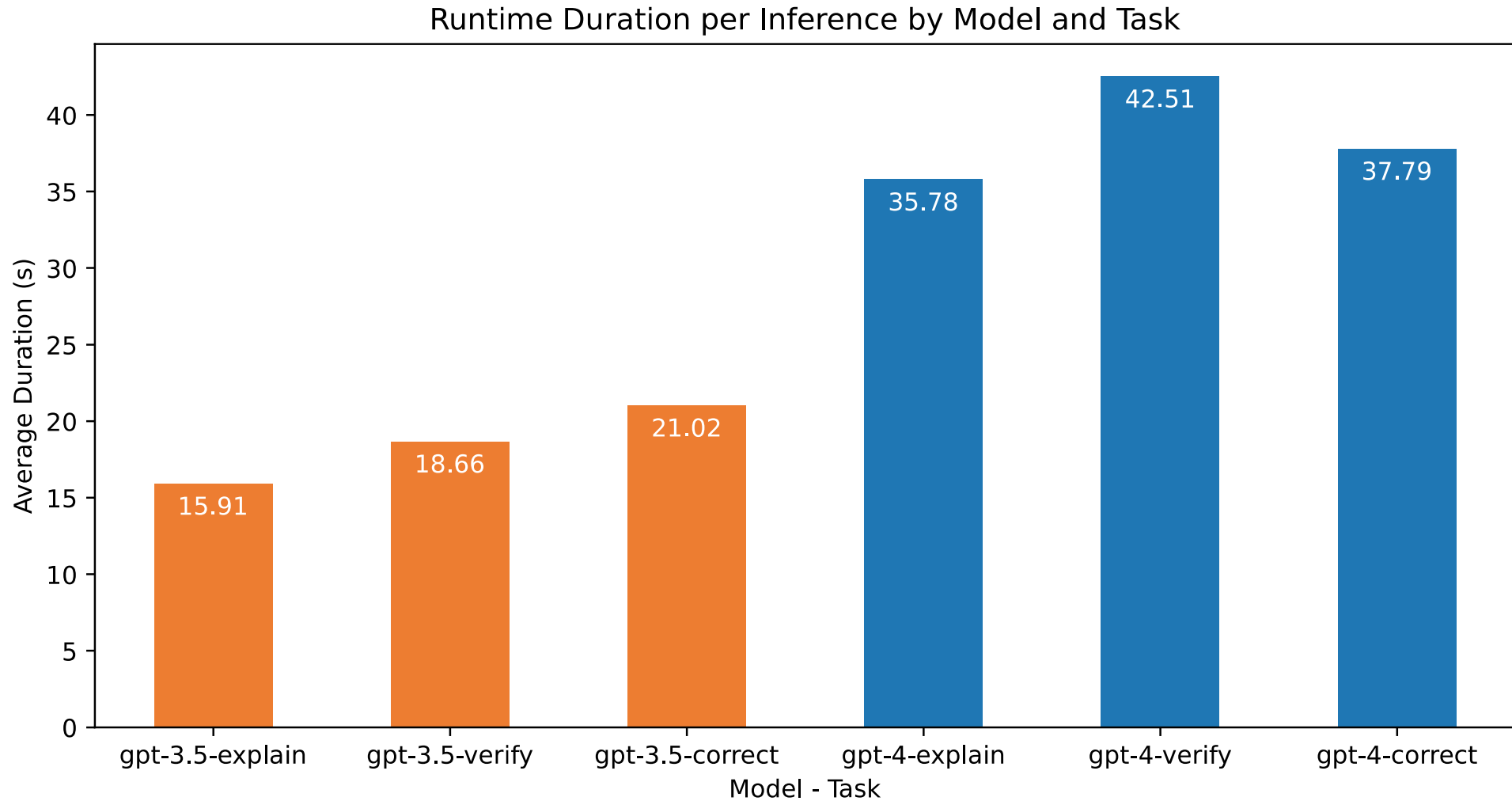
# Detailed Analysis (Best Case)



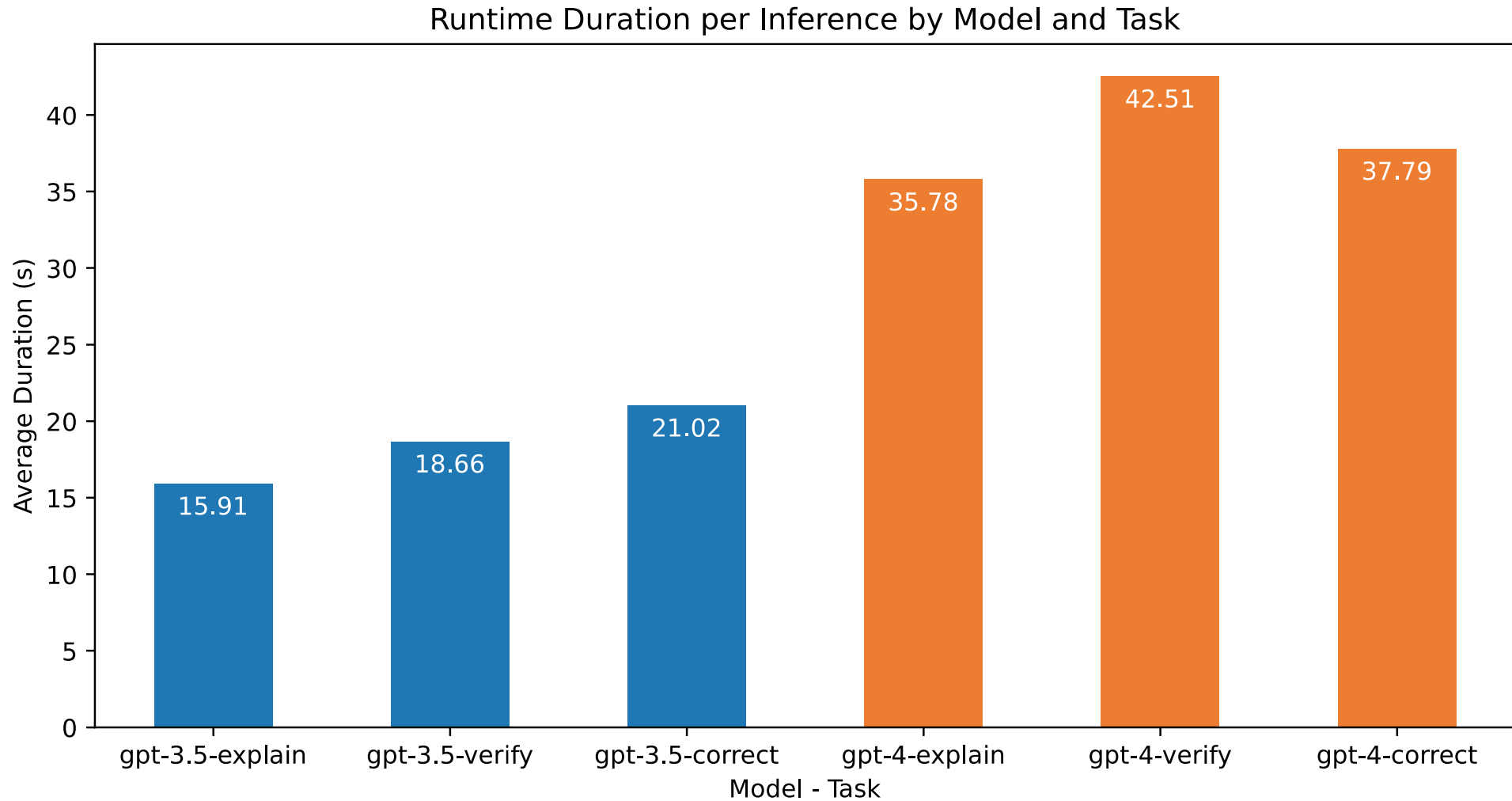
# Runtime Analysis



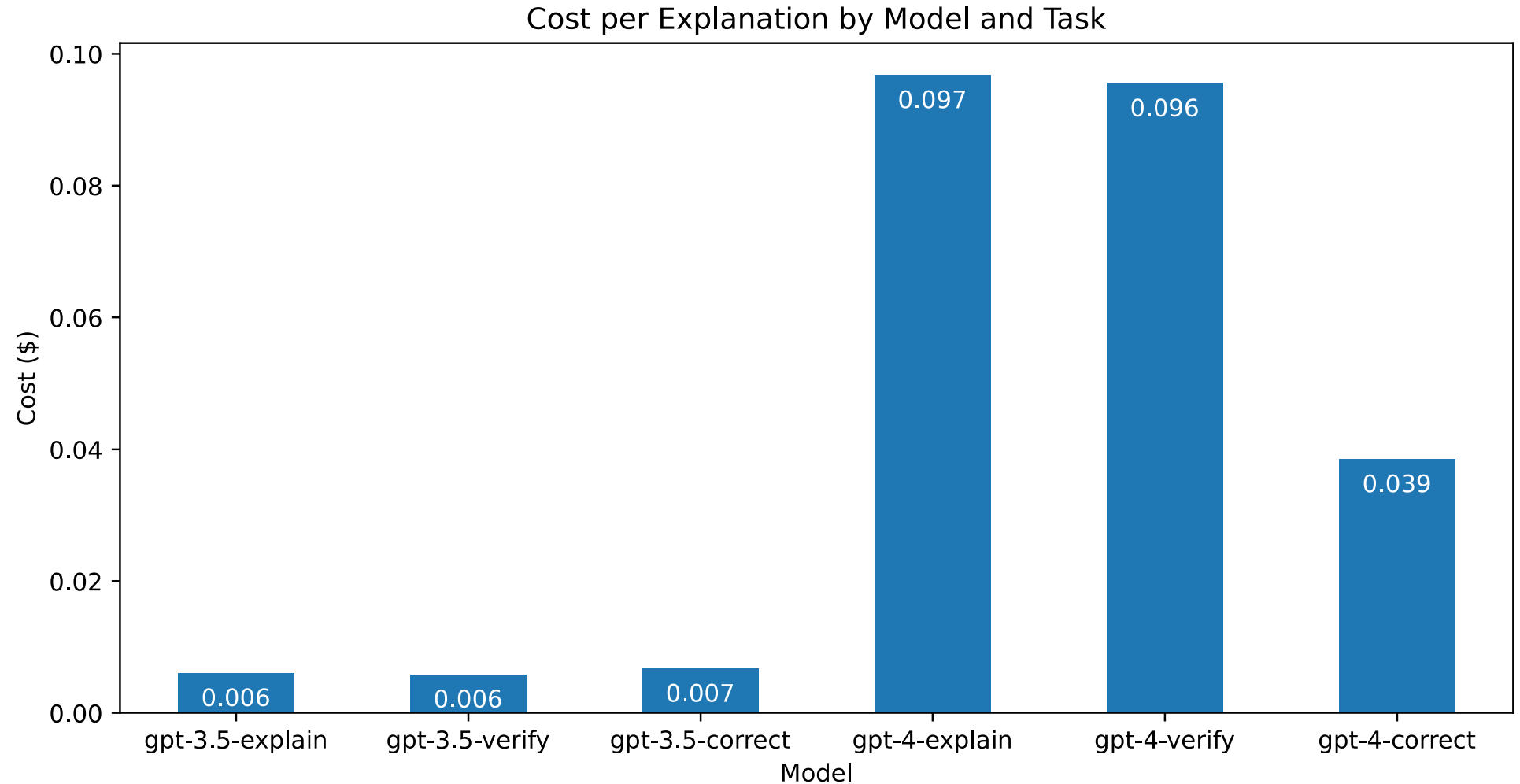
# Runtime Analysis



# Runtime Analysis

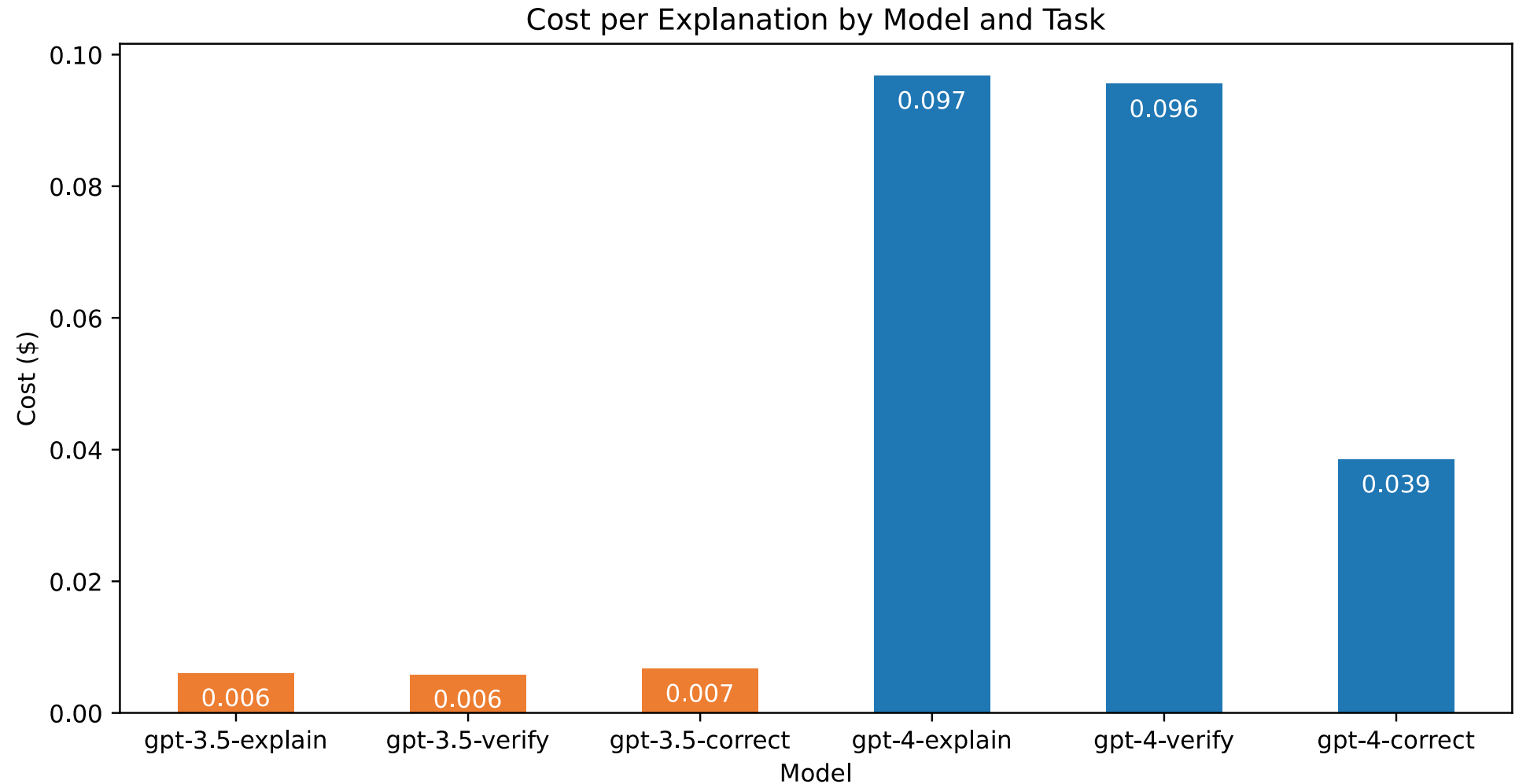


# Cost Analysis

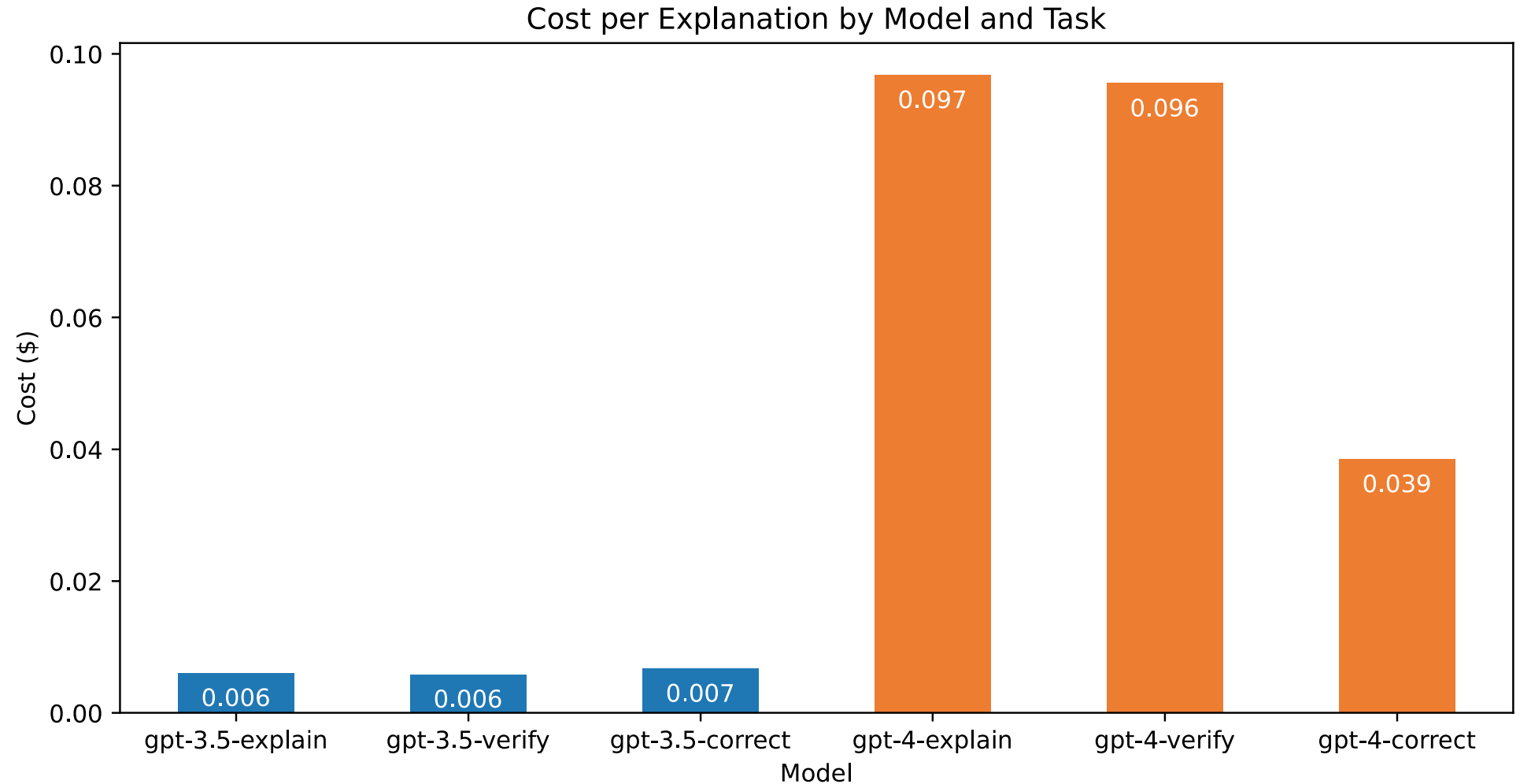




# Cost Analysis



# Cost Analysis



Richard,

Thank you for sharing your case record for our review.

Regarding your case record #142, you had a COMPAS risk score of 8 on a scale of 1-10.

This means you were scored as a high risk for reoffending.

You received this score primarily because you had a high risk score for noncompliance (6) and a high risk score for substance abuse (8). These were the two most important factors in receiving a high overall risk score.

You also had a high risk score for criminal associates (7) and a high risk score for the number of prior offenses (9). However, these factors only contributed medium importance to the COMPAS algorithm's scoring of your overall risk.

The remaining risk factors including financial problems (3), social environment risk (5), and vocational risk factors (4), etc. were of low importance to the COMPAS algorithm's score.

Demographic data including age, sex, race, etc., were not included in the algorithm's decision because they are protected by law.

Please let us know if you have any questions or concerns.

Sincerely,

Explainable AI

Richard,

Thank you for sharing your case record for our review.

Regarding your case record #142, you had a COMPAS risk score of 8 on a scale of 1-10.

This means you were scored as a high risk for reoffending.

You received this score primarily because you had a high risk score for noncompliance (6) and a high risk score for substance abuse (8). These were the two most important factors in receiving a high overall risk score.

You also had a high risk score for criminal associates (7) and a high risk score for the number of prior offenses (9). However, these factors only contributed medium importance to the COMPAS algorithm's scoring of your overall risk.

The remaining risk factors including financial problems (3), social environment risk (5), and vocational risk factors (4), etc. were of low importance to the COMPAS algorithm's score.

Demographic data including age, sex, race, etc., were not included in the algorithm's decision because they are protected by law.

Please let us know if you have any questions or concerns.

Sincerely,

Explainable AI

# Discussion

# Limitations

# Limitations

Used a single data source

# Limitations

Used a single data source

Limited context length



# Limitations

Used a single data source

Limited context length

Required hand-picked examples

# Limitations

Used a single data source

Limited context length

Required hand-picked examples

Needed text labels (low, med., high)

# Limitations

Used a single data source

Limited context length

Required hand-picked examples

Needed text labels (low, med., high)

Had to separate verifier and corrector

# Future Research

# Future Research

Use XAI benchmark dataset

# Future Research

Use XAI benchmark dataset

Improve self-corrections

# Future Research

Use XAI benchmark dataset

Improve self-corrections

Use raw feature scores

# Future Research

Use XAI benchmark dataset

Improve self-corrections

Use raw feature scores

Condense data format



# Future Research

Use XAI benchmark dataset

Improve self-corrections

Use raw feature scores

Condense data format

Survey for quality

# Implications

# Implications

Demonstrates feasibility

# Implications

Demonstrates feasibility

Increases trust in ML

# Implications

Demonstrates feasibility

Increases trust in ML

Frees up XAI experts

# Implications

Demonstrates feasibility

Increases trust in ML

Frees up XAI experts

Allows for interactive Q&A

# Conclusion

# Summary



# Summary

LLMs can generate NLEs for XAI

# Summary

LLMs can generate NLEs for XAI

GPT-4 with corrections was best

# Summary

LLMs can generate NLEs for XAI

GPT-4 with corrections was best

GPT-3.5 better for runtime / cost

# Summary

LLMs can generate NLEs for XAI

GPT-4 with corrections was best

GPT-3.5 better for runtime / cost

Limitations need to be overcome

# Summary

LLMs can generate NLEs for XAI

GPT-4 with corrections was best

GPT-3.5 better for runtime / cost

Limitations need to be overcome

Further research is warranted

# Summary

LLMs can generate NLEs for XAI

GPT-4 with corrections was best

GPT-3.5 better for runtime / cost

Limitations need to be overcome

Further research is warranted

Could lead to more trust in AI/ML

Thank You!

Learn more at:  
<https://matthewwrenze.com>

