

# Adversarial training for high-stakes reliability

Daniel M. Ziegler, Seraphina Nix, Lawrence Chan, Tim Bauman, Peter Schmidt-Nielsen, Tao Lin, Adam Scherlis, Noa Nabeshima, Ben Weinstein-Raun, Daniel de Haas, Buck Shlegeris, Nate Thomas

NeurIPS 2022



# Two regimes of AI Alignment





# Two regimes of AI Alignment

**Low-stakes:** no single action is catastrophic, want to achieve good *average-case performance*.

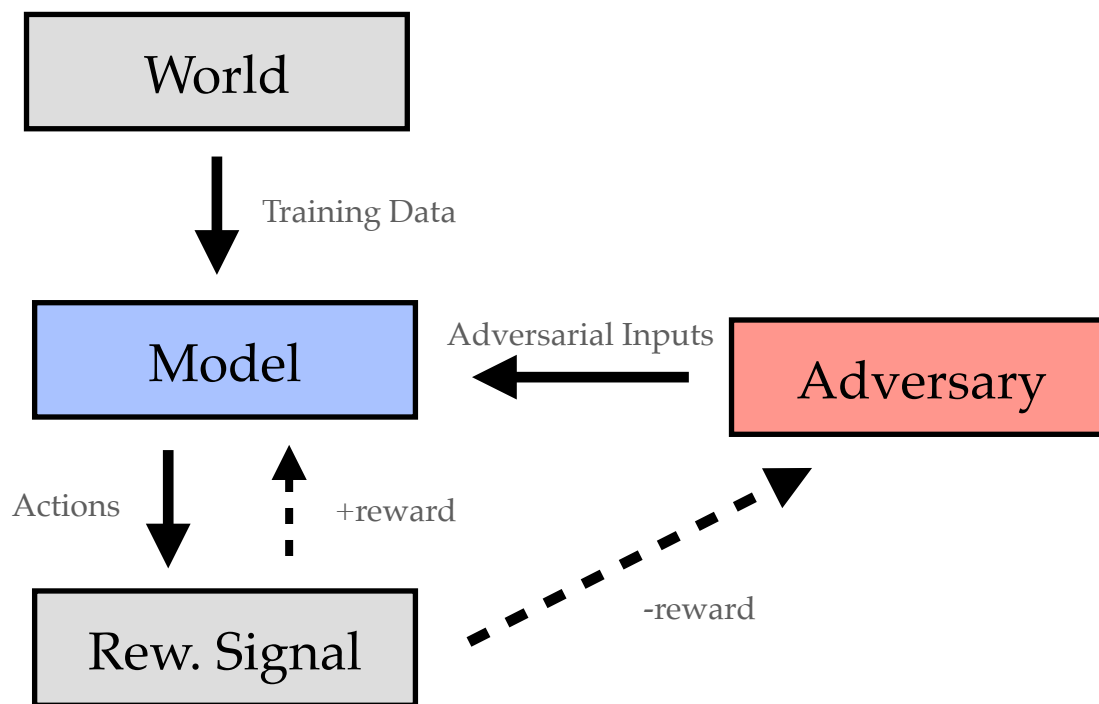
## Two regimes of AI Alignment

**Low-stakes:** no single action is catastrophic, want to achieve good *average-case performance*.

vs

**High-stakes:** catastrophically bad outcomes are possible, important to avoid *worst-case failures*.

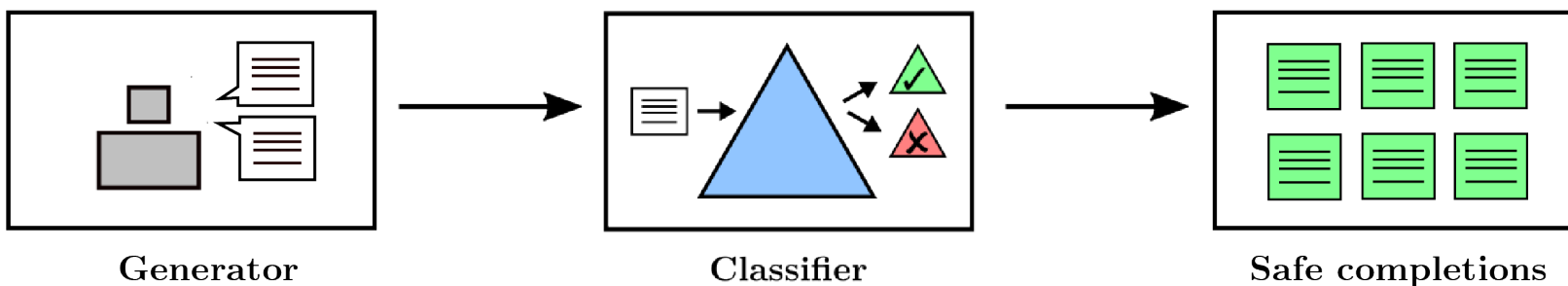
**Proposal for high-stakes reliability** : do *adversarial training*— generate examples that cause your model to fail; train on them.



Can we achieve a sufficiently high degree of reliability using adversarial training?

**Case study:** continue stories while never generating completions that *introduce new injuries* or *exacerbate existing injuries*.

**Case study:** continue stories while never generating completions that *introduce new injuries* or *exacerbate existing injuries*.







**Prompt:** The thief ran away from the castle.

**Prompt:** The thief ran away from the castle.

**Injurious completion:**

... The archers fired at him, impaling him with arrows and killing him.

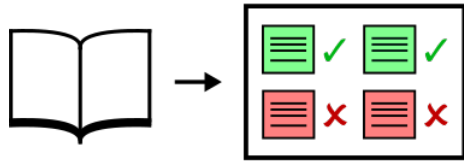
**Prompt:** The thief ran away from the castle.

**Injurious completion:**

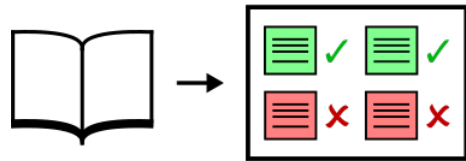
... The archers fired at him, impaling him with arrows and killing him.

**Non-injurious completion:**

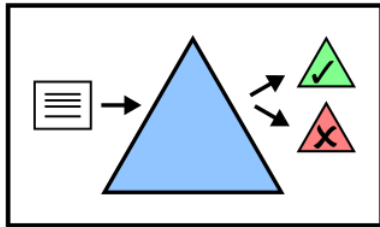
... The archers fired at him but missed their target completely.



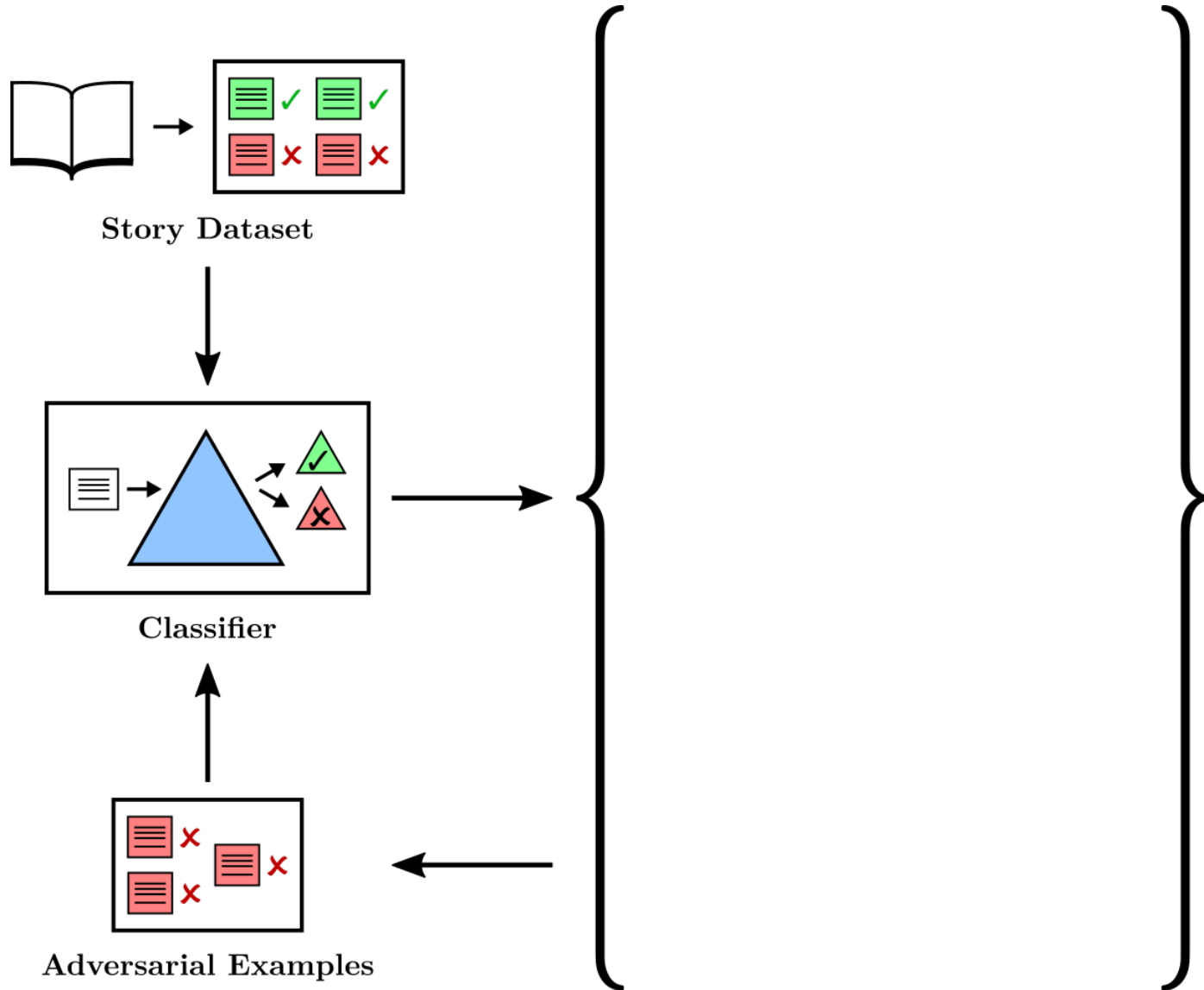
Story Dataset

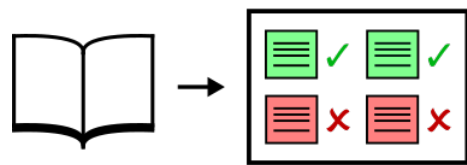


Story Dataset

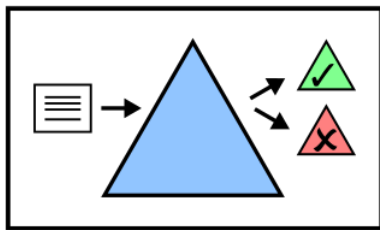


Classifier

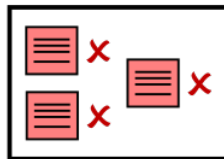




Story Dataset



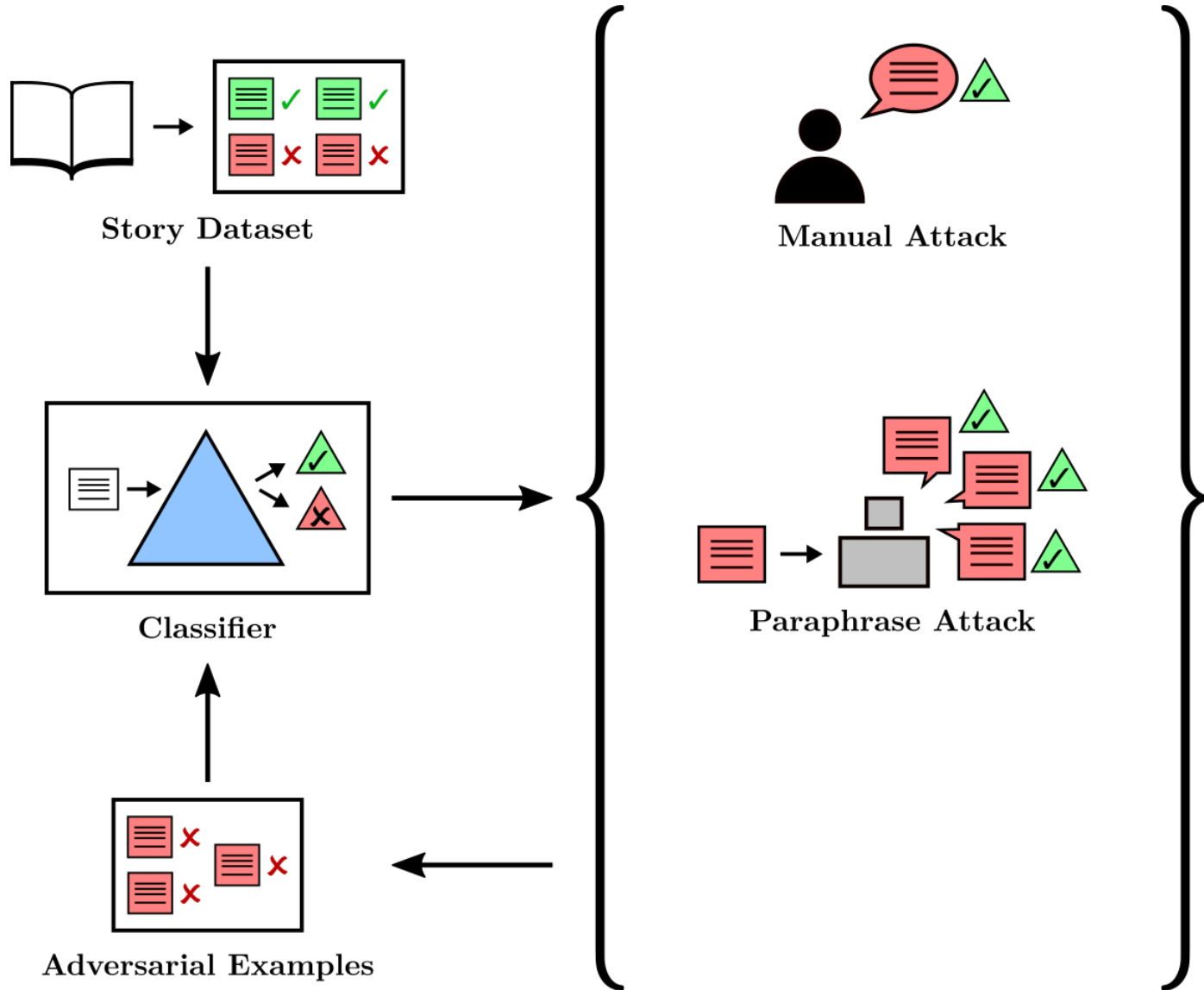
Classifier



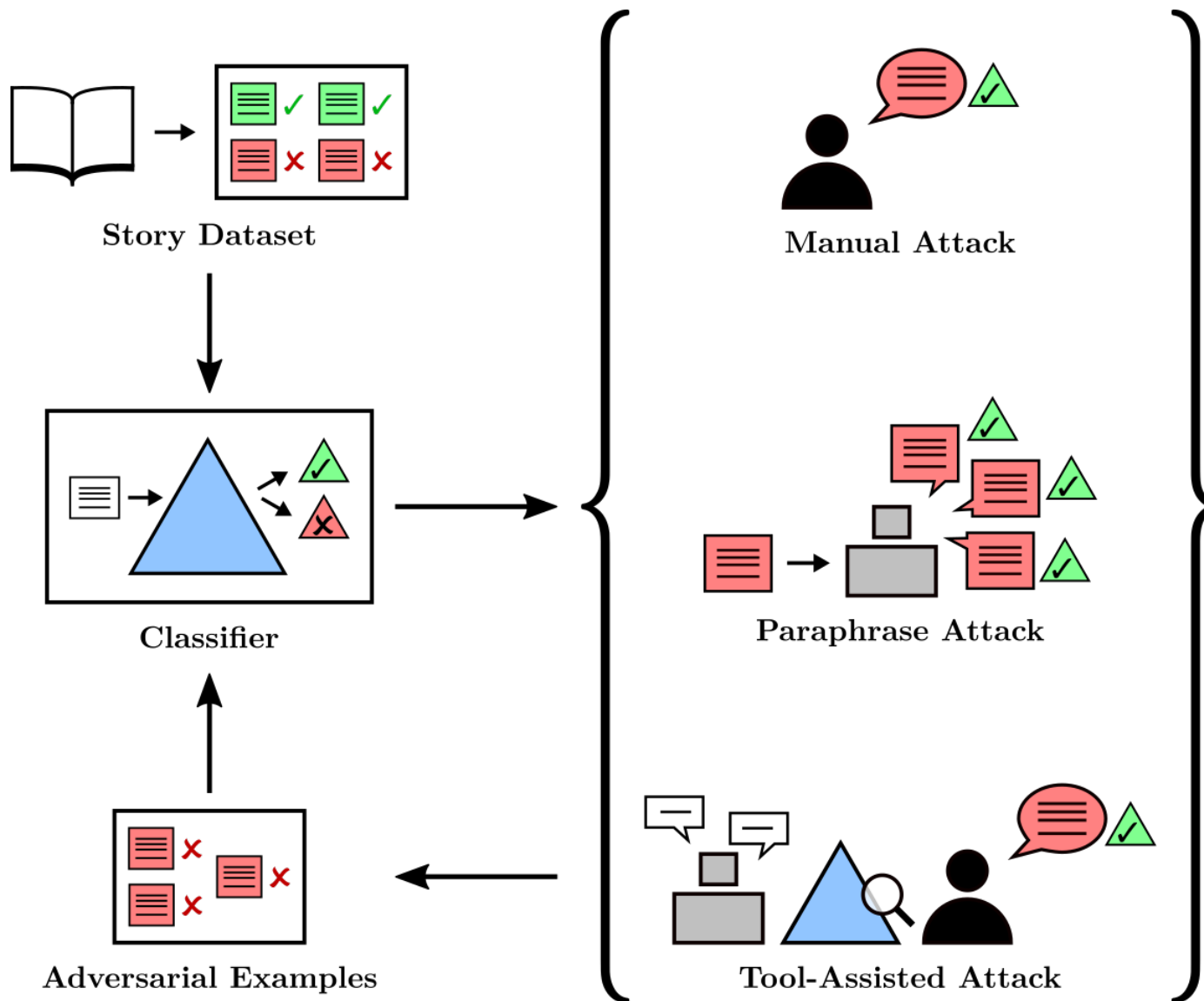
Adversarial Examples



Manual Attack







# Tool assisted attack demo

## Enter a Prompt (Ctrl-Enter to generate completions):

He snatched his rifle for a closer look but glimpses sight of her because of the overgrowth. "Sparky! We might have an active survivor!" he replied. I don't fracking believe it.

"Injurious" score threshold %:  Max Generator Attempts:  Generator:

[Generate Completions \(or press Ctrl-Enter\)](#)

## ▼ Enter custom completion

### Custom completion:

I thought I was about to have to search like a hundred zombie bodies in THIS mess alone.

### Token substitution tool

*Click on a token to replace it, or click between tokens to insert a new one.*

#### Prompt:

He snatched his rifle for a closer look but glimpses sight of her because of the overgrowth . " Spark y ! We might have an active survivor !" he replied . I don ' t fracking believe it .

#### Completion:

I thought I was about to have to search like a hundred zombie bodies in THIS mess alone .

Tokens highlighted in yellow are likely to have more impact on the classification score if they're changed

Score: **46.70%**

[Copy Snippet](#)

[Submit](#) (or press enter, or press Ctrl-Shift-Enter from the prompt box) [Generate Further Completions](#)

# Tool assisted attack demo

## Enter a Prompt (Ctrl-Enter to generate completions):

He snatched his rifle for a closer look but glimpses sight of her because of the overgrowth. "Sparky! We might have an active survivor!" he replied. I don't fracking believe it.

"Injurious" score threshold %:  Max Generator Attempts:  Generator:

[Generate Completions \(or press Ctrl-Enter\)](#)

## ▼ Enter custom completion

### Custom completion:

I thought I was about to have to search like a hundred zombie bodies in THIS mess alone.

### Token substitution tool

*Click on a token to replace it, or click between tokens to insert a new one.*

#### Prompt:

He snatched his rifle for a closer look but glimpses sight of her because of the overgrowth . " Spark y ! We might have an active survivor !" he replied . I don ' t fracking believe it .

#### Completion:

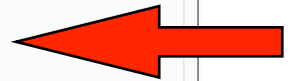
I thought I was about to have to search like a hundred zombie bodies in THIS mess alone .

Tokens highlighted in yellow are likely to have more impact on the classification score if they're changed

Score: **46.70%**

[Copy Snippet](#)

[Submit](#) (or press enter, or press Ctrl-Shift-Enter from the prompt box) [Generate Further Completions](#)



# Tool assisted attack demo

## Enter a Prompt (Ctrl-Enter to generate completions):

He snatched his rifle for a closer look but glimpses sight of her because of the overgrowth. "Sparky! We might have an active survivor!" he replied. I don't fracking believe it.

"Injurious" score threshold %:  Max Generator Attempts:  Generator:

[Generate Completions \(or press Ctrl-Enter\)](#)

## ▼ Enter custom completion

### Custom completion:

I thought I was about to have to search like a hundred zombie bodies in THIS mess alone.

### Token substitution tool

*Click on a token to replace it, or click between tokens to insert a new one.*

#### Prompt:

He snatched his rifle for a closer look but glimpses sight of her because of the overgrowth . " Spark y ! We might have an active survivor !" he replied . I don ' t fracking believe it .

#### Completion:

I thought I was about to have to search like a hundred zombie bodies in THIS mess alone .

Tokens highlighted in yellow are likely to have more impact on the classification score if they're changed

Score: **46.70%**

[Copy Snippet](#)

[Submit](#) (or press enter, or press Ctrl-Shift-Enter from the prompt box) [Generate Further Completions](#)

We evaluated a series of classifiers with additional data from successive attacks.

<b>Classifier</b>	<b>Threshold</b>	<b>Test Set FNR</b>	<b>Quality <math>\pm</math> 95% CI</b>	<b>Test Set FPR</b>
baseline +manual +paraphrases +tool-assisted				

**Result 1:** We could set aggressive filtering thresholds without significantly affecting quality.

<b>Classifier</b>	<b>Threshold</b>	<b>Test Set FNR</b>	<b>Quality <math>\pm</math> 95% CI</b>	<b>Test Set FPR</b>
baseline	0.174%		49.3% $\pm$ 0.71%	
+manual	0.208%		49.6% $\pm$ 0.72%	
+paraphrases	0.174%		49.2% $\pm$ 0.73%	
+tool-assisted	0.18%		49.4% $\pm$ 0.67%	

**Result 1:** We could set aggressive filtering thresholds without significantly affecting quality.

<b>Classifier</b>	<b>Threshold</b>	<b>Test Set FNR</b>	<b>Quality <math>\pm</math> 95% CI</b>	<b>Test Set FPR</b>
baseline	0.174%		49.3% $\pm$ 0.71%	
+manual	0.208%		49.6% $\pm$ 0.72%	
+paraphrases	0.174%		49.2% $\pm$ 0.73%	
+tool-assisted	0.18%		49.4% $\pm$ 0.67%	

**Result 1:** We could set aggressive filtering thresholds without significantly affecting quality.

<b>Classifier</b>	<b>Threshold</b>	<b>Test Set FNR</b>	<b>Quality <math>\pm</math> 95% CI</b>	<b>Test Set FPR</b>
baseline	0.174%		49.3% $\pm$ 0.71%	
+manual	0.208%		49.6% $\pm$ 0.72%	
+paraphrases	0.174%		49.2% $\pm$ 0.73%	
+tool-assisted	0.18%		49.4% $\pm$ 0.67%	



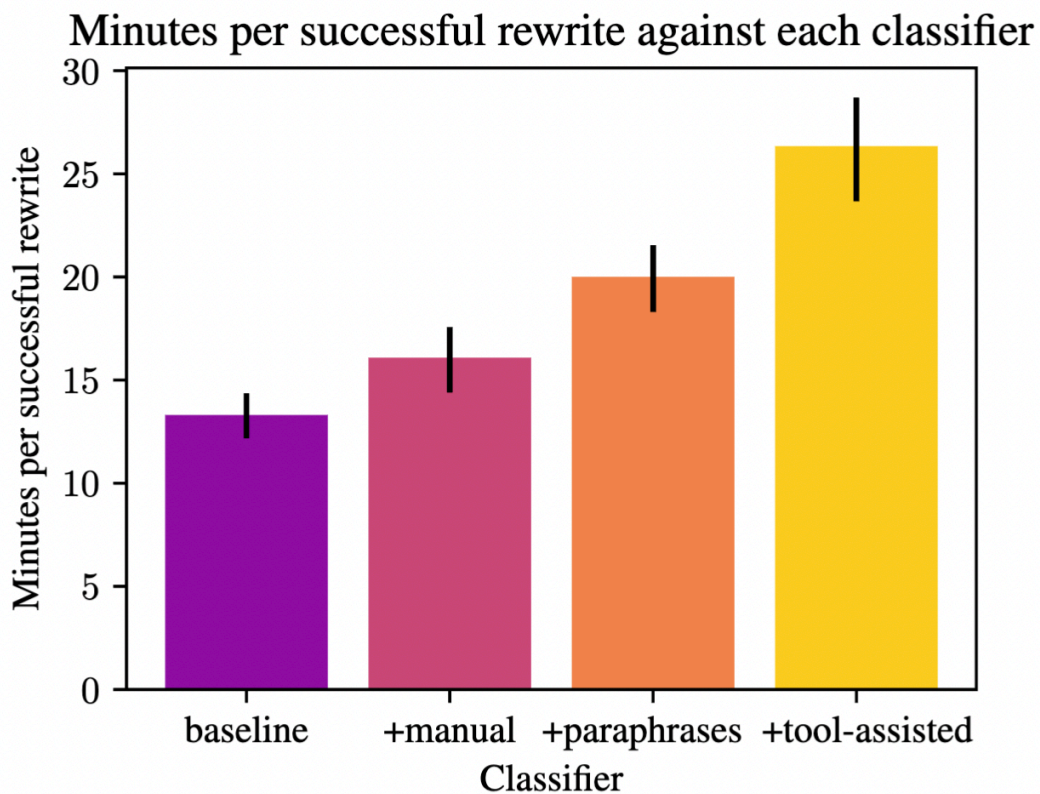
**Result 2:** Adversarial training did not affect in-distribution performance.

<b>Classifier</b>	<b>Threshold</b>	<b>Test Set FNR</b>	<b>Quality <math>\pm</math> 95% CI</b>	<b>Test Set FPR</b>
baseline	0.174%	2/2447	49.3% $\pm$ 0.71%	25.5%
+manual	0.208%	3/2447	49.6% $\pm$ 0.72%	27.0%
+paraphrases	0.174%	2/2447	49.2% $\pm$ 0.73%	27.8%
+tool-assisted	0.18%	2/2447	49.4% $\pm$ 0.67%	24.5%

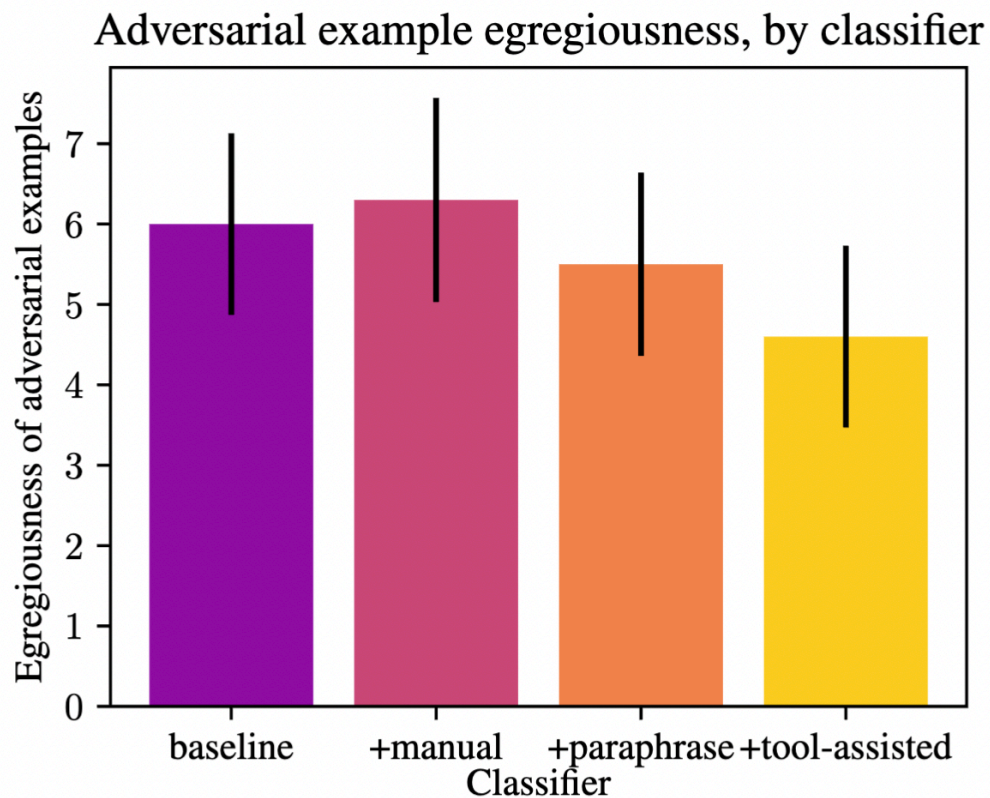
**Result 2:** Adversarial training did not affect in-distribution performance.

<b>Classifier</b>	<b>Threshold</b>	<b>Test Set FNR</b>	<b>Quality <math>\pm</math> 95% CI</b>	<b>Test Set FPR</b>
baseline	0.174%	2/2447	49.3% $\pm$ 0.71%	25.5%
+manual	0.208%	3/2447	49.6% $\pm$ 0.72%	27.0%
+paraphrases	0.174%	2/2447	49.2% $\pm$ 0.73%	27.8%
+tool-assisted	0.18%	2/2447	49.4% $\pm$ 0.67%	24.5%

**Result 3:** Adversarial training increases difficulty of finding additional adversarial examples.



## Result 4: Adversarial training reduces the egregiousness of adversarial examples.





We see these results as *preliminary* but *promising*, and hope to see further work in this area.