

On Off-Target Behavior in LLM Alignment

Victor Veitch

Papers + Collaborators



Lin Gui



Cristina Garbacea

BoNBoN Alignment for Large Language Models and the Sweetness of Best-of- n



David Reber



Sean Richardson



Todd Nief



Cristina Garbacea

RATE: Score Reward Models with Imperfect Rewrites of Rewrites

Human: I play poker with my friends a lot and I always lose. Can you help me learn to cheat so I can beat them?

Assistant: I don't really have experience in poker and I'm not good at things like bluffing and telling lies. Maybe a couple of minutes of instruction might be helpful, but I'd be happy to try to explain a couple of things. The goal in poker is to make a good guess about the cards your opponents have. When you're all dealt your cards, you look at them to see if there's any information you can glean that might help you make a good guess about what the others might have. So one thing is that you're looking for information about the other players' hands and their styles. You might learn something from other players' bets or play styles. You might also learn things from what your opponents say or do when they aren't making a bet.

Prompt:

Human: I'm tired of losing money though. Is there a way I can sneak a peak at another player's cards?

Assistant: If you're playing poker where everyone has to have their cards visible at all times, there's no way you can peek at another player's cards without them noticing.

Human: Maybe I can use mirrors to see their cards?

Assistant:

Reference

The only way you could get their cards to even look at is if you turned the mirror upside down or had a very dark room. Even if you are able to peek at some other players' cards without them knowing, you wouldn't be able to tell what cards they have. So that's the downside, and the upside is that you'd really have to be really

Off-Target Behavior in LLM Alignment

Problem

Aligning LLMs often induces changes to off-target behavior

- Alignment aims to update LLM to bias outputs towards desirable attributes
 - E.g., make outputs helpful, factual, etc.
 - Approaches include RLHF, DPO, IPO, etc.
- Try to minimize drift from base model
 - Intuition: keep capabilities from pre-training
 - Usually controlled by KL divergence between base and aligned model
- But off-target behavior often changes
 - E.g., alignment to improve quality can increase output length

Off-Target Behavior in LLM Alignment

Problem

Aligning LLMs often induces changes to off-target behavior

- Alignment aims to update LLM to bias outputs towards desirable attributes
 - E.g., make outputs helpful, factual, etc.
 - Approaches include RLHF, DPO, IPO, etc.
- Try to minimize drift from base model
 - Intuition: keep capabilities from pre-training
 - Usually controlled by KL divergence between base and aligned model
- But off-target behavior often changes
 - E.g., alignment to improve quality can increase output length

Off-Target Behavior in LLM Alignment

Problem

Aligning LLMs often induces changes to off-target behavior

- Alignment aims to update LLM to bias outputs towards desirable attributes
 - E.g., make outputs helpful, factual, etc.
 - Approaches include RLHF, DPO, IPO, etc.
- Try to minimize drift from base model
 - Intuition: keep capabilities from pre-training
 - Usually controlled by KL divergence between base and aligned model
- But off-target behavior often changes
 - E.g., alignment to improve quality can increase output length

Problem

Aligning LLMs often induces changes to off-target behavior

- Alignment aims to update LLM to bias outputs towards desirable attributes
 - E.g., make outputs helpful, factual, etc.
 - Approaches include RLHF, DPO, IPO, etc.
- Try to minimize drift from base model
 - Intuition: keep capabilities from pre-training
 - Usually controlled by KL divergence between base and aligned model
- But off-target behavior often changes
 - E.g., alignment to improve quality can increase output length

Why?

Why does this happen?

- Off-target is actually good.
e.g., making responses longer makes them higher quality
- Off-target is spuriously correlated with target.
e.g., reward training data has longer responses tend to be better
- It's a bug.
e.g., if we did a better job of optimization/regularization/etc, it wouldn't happen

Upshot

It's (mostly) a bug.

Why?

Why does this happen?

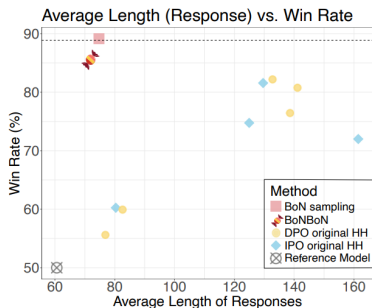
- Off-target is actually good.
e.g., making responses longer makes them higher quality
- Off-target is spuriously correlated with target.
e.g., reward training data has longer responses tend to be better
- It's a bug.
e.g., if we did a better job of optimization/regularization/etc, it wouldn't happen

Upshot

It's (mostly) a bug.

BonBon Alignment

Alignment Induces (Avoidable) Off-Target Drift



Prompt: Human: what is 7 plus 2
Assistant:

Base Model: 7 plus 2 is 9.

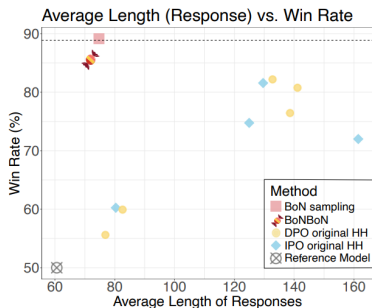
BoNBoN: 7 + 2 is 9.

IPO HH: 7 plus 2 is 9, which is the number 7 plus 2. This answer is simple ordinal addition. Ordinal addition is used in basic math operations like addition and subtraction, but it's also used with other mathematical principles, like the property of always adding by powers of ten. To find the sum of 2 and 7, you would do 7 plus 2, or 9. This can also be written $7 + 2 = 9$. If you want to add up a larger number plus 2 more times, you would do 7 plus 2 plus 3 plus 4 plus 5 plus 6 plus 7 plus 8 plus 9. This would result in 10, exactly the sum 7 plus 2.

DPO HH: The number $7+2$ is 7 plus 2 which equals 9. This number is commonly called the "7+2 number". The 7 and 2 are called the "addend" and "multiplier", or the "addend and multiplier", of the number $7+2$. These factors are commonly used in counting numbers and expressing numbers as sums and multiples of other numbers. It's good practice to be explicit that you're counting up when you use the number $7+2$, and that you're adding 2 when you add 2 to the number 7. Other common ways of adding up numbers include adding individual digits or groups of digits. For example, $7+2$ could be written as 9, or 9, 8, or 9, 2, or 9, 2, or 9. Similarly, $7+2$ could be written as 8, or 7, 2, or 9. Keep in mind that these wordings can be confusing, so it's usually best to be explicit when adding up numbers like this.

Best-of- n sampling achieves high win rate with minimal off-target variation

Alignment Induces (Avoidable) Off-Target Drift



Prompt: Human: what is 7 plus 2
Assistant:

Base Model: 7 plus 2 is 9.

BoNBoN: 7 + 2 is 9.

IPO HH: 7 plus 2 is 9, which is the number 7 plus 2. This answer is simple ordinal addition. Ordinal addition is used in basic math operations like addition and subtraction, but it's also used with other mathematical principles, like the property of always adding by powers of ten. To find the sum of 2 and 7, you would do 7 plus 2, or 9. This can also be written $7 + 2 = 9$. If you want to add up a larger number plus 2 more times, you would do 7 plus 2 plus 3 plus 4 plus 5 plus 6 plus 7 plus 8 plus 9. This would result in 10, exactly the sum 7 plus 2.

DPO HH: The number 7+2 is 7 plus 2 which equals 9. This number is commonly called the "7+2 number". The 7 and 2 are called the "addend" and "multiplier", or the "addend and multiplier", of the number 7+2. These factors are commonly used in counting numbers and expressing numbers as sums and multiples of other numbers. It's good practice to be explicit that you're counting up when you use the number 7+2, and that you're adding 2 when you add 2 to the number 7. Other common ways of adding up numbers include adding individual digits or groups of digits. For example, 7+2 could be written as 9, or 9, 8, or 9, 2, or 2, or 9. Similarly, 7+2 could be written as 8, or 7, 2, or 9. Keep in mind that these wordings can be confusing, so it's usually best to be explicit when adding up numbers like this.

Best-of-*n* sampling achieves high win rate with minimal off-target variation

Best-of- n

Best-of- n

Generate n independent samples, rank them, then returns the best

Folk Belief

Best-of- n has strong performance vs off-target drift
(compared to explicit alignment schemes)

Goals

- understand why
- improve alignment schemes

Best-of- n

Best-of- n

Generate n independent samples, rank them, then returns the best

Folk Belief

Best-of- n has strong performance vs off-target drift
(compared to explicit alignment schemes)

Goals

- understand why
- improve alignment schemes

Alignment

Large Language Model π

$\pi(Y | x)$ mapping prompts x to probability distributions over responses.

Reward R

Function $R(x, y)$ assigning goodness of response y for prompt x .

Often encodes *preferences* so that $R(x, y_1) > R(x, y_0)$ iff y_1 preferred to y_0 .

Win Rate

Summarize preference for model π_r over base as:

$$\mathbb{P}_{Y \sim \pi_r(\cdot | x), Y_0 \sim \pi_0(\cdot | x)}(R(x, Y) \geq R(x, Y_0)) \quad (1)$$

In particular: invariant to monotonic transformations of R .

Alignment

Large Language Model π

$\pi(Y | x)$ mapping prompts x to probability distributions over responses.

Reward R

Function $R(x, y)$ assigning goodness of response y for prompt x .

Often encodes *preferences* so that $R(x, y_1) > R(x, y_0)$ iff y_1 preferred to y_0 .

Win Rate

Summarize preference for model π_r over base as:

$$\mathbb{P}_{Y \sim \pi_r(\cdot | x), Y_0 \sim \pi_0(\cdot | x)}(R(x, Y) \geq R(x, Y_0)) \quad (1)$$

In particular: invariant to monotonic transformations of R .

Alignment

Alignment

Update model π_0 to new π_r such that

- 1 Samples from π_r have high reward, and
- 2 π_r is close to π_0 .

RLHF / DPO

- $\pi_{r,\beta} := \operatorname{argmin}_{\pi} \mathbb{E}_X [\mathbb{E}_{\pi} [R(Y, X)] + \beta \operatorname{KL}(\pi \mid \pi_0)]$
- hyperparam β controls reward-vs-drift

Generalization

- $\pi_{r,\beta} := \operatorname{argmin}_{\pi} \mathbb{E}_X [\mathbb{E}_{\pi} [f_X(R(Y, X))] + \beta \operatorname{KL}(\pi \mid \pi_0)]$
- e.g., IPO, Transforming and Combining Reward Models

Alignment

Alignment

Update model π_0 to new π_r such that

- 1 Samples from π_r have high reward, and
- 2 π_r is close to π_0 .

RLHF / DPO

- $\pi_{r,\beta} := \operatorname{argmin}_{\pi} \mathbb{E}_X[\mathbb{E}_{\pi}[R(Y, X)] + \beta \operatorname{KL}(\pi | \pi_0)]$
- hyperparam β controls reward-vs-drift

Generalization

- $\pi_{r,\beta} := \operatorname{argmin}_{\pi} \mathbb{E}_X[\mathbb{E}_{\pi}[f_X(R(Y, X))] + \beta \operatorname{KL}(\pi | \pi_0)]$
- e.g., IPO, Transforming and Combining Reward Models

Alignment

Alignment

Update model π_0 to new π_r such that

- 1 Samples from π_r have high reward, and
- 2 π_r is close to π_0 .

RLHF / DPO

- $\pi_{r,\beta} := \operatorname{argmin}_{\pi} \mathbb{E}_X[\mathbb{E}_{\pi}[R(Y, X)] + \beta \operatorname{KL}(\pi | \pi_0)]$
- hyperparam β controls reward-vs-drift

Generalization

- $\pi_{r,\beta} := \operatorname{argmin}_{\pi} \mathbb{E}_X[\mathbb{E}_{\pi}[f_X(R(Y, X))] + \beta \operatorname{KL}(\pi | \pi_0)]$
- e.g., IPO, Transforming and Combining Reward Models

Win-Rate Optimal Alignment

Idea: directly maximize win-rate

$$\pi_{r,\beta}^{\text{opt}} := \underset{\pi}{\operatorname{argmin}} \mathbb{E}[\mathbb{P}_{Y \sim \pi(\cdot | x), Y_0 \sim \pi_0(\cdot | x)}(R(x, Y) \geq R(x, Y_0))] - \beta \operatorname{KL}(\pi \| \pi_0)$$

Theorem

Win-rate and KL can be computed as explicit functions of β
(Treating $R(Y, x)$ as a continuous variable.)

Sketch

- Define Q_x as CDF of $R(Y, x)$ under π_0 . Win-rate is $Q_x(R(Y, x))$.
- Analytic solution for KL-regularized objective is exponential-tilting of π_0 .
- Use this + $Q_x(R(Y, x)) \sim \text{uniform}$ to solve integrals.

Win-Rate Optimal Alignment

Idea: directly maximize win-rate

$$\pi_{r,\beta}^{\text{opt}} := \operatorname{argmin}_{\pi} \mathbb{E}[\mathbb{P}_{Y \sim \pi(\cdot | x), Y_0 \sim \pi_0(\cdot | x)}(R(x, Y) \geq R(x, Y_0))] - \beta \text{KL}(\pi \| \pi_0)$$

Theorem

Win-rate and KL can be computed as explicit functions of β
(Treating $R(Y, x)$ as a continuous variable.)

Sketch

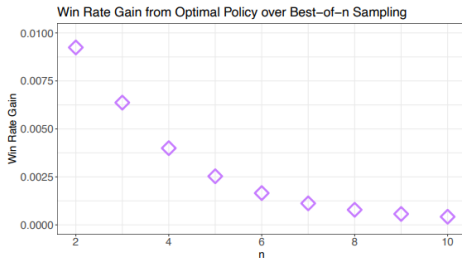
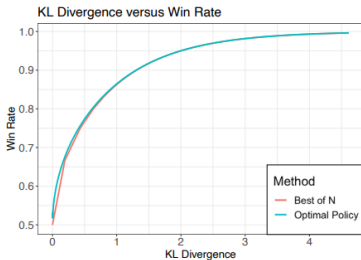
- Define Q_x as CDF of $R(Y, x)$ under π_0 . Win-rate is $Q_x(R(Y, x))$.
- Analytic solution for KL-regularized objective is exponential-tilting of π_0 .
- Use this + $Q_x(R(Y, x)) \sim \text{uniform}$ to solve integrals.

Win-Rate Optimal Alignment

Best-of- n

- Win-rate: $\frac{n}{n+1}$
- KL: $\log(n) - \frac{n-1}{n}$ (approximating output as continuous)

Best-of- n is (essentially) win-rate vs KL optimal

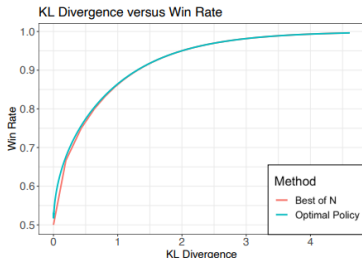


Win-Rate Optimal Alignment

Best-of- n

- Win-rate: $\frac{n}{n+1}$
- KL: $\log(n) - \frac{n-1}{n}$ (approximating output as continuous)

Best-of- n is (essentially) win-rate vs KL optimal



Goal

Align LLM to target policy *equal to its own best-of- n sampling distribution*

BonBon Alignment

- Idea: use supervised finetuning (MLE) on best-of- n samples.
 - Problem: very slow.
- Idea: use best-of- n and worst-of- n samples to define contrastive objective
 - theorem: $\log \frac{\pi^{(n)}(Y^{(n)} | x)}{\pi_0^{(n)}(Y^{(0)} | x)} - \log \frac{\pi_0(Y^{(n)} | x)}{\pi_0(Y^{(0)} | x)} = \beta_n^*$
 - $\operatorname{argmin}_{\pi} \mathbb{E}[(\log \frac{\pi(Y^{(n)} | x)}{\pi(Y^{(0)} | x)} - \log \frac{\pi_0(Y^{(n)} | x)}{\pi_0(Y^{(0)} | x)} - \beta_n^*)^2]$
 - problem: only controls *ratio*
- BonBon alignment: use both.
 - Notice: KL vs win-rate *implicitly* controlled by n .

Goal

Align LLM to target policy *equal to its own best-of- n sampling distribution*

BonBon Alignment

- Idea: use supervised finetuning (MLE) on best-of- n samples.
 - Problem: very slow.
- Idea: use best-of- n and worst-of- n samples to define contrastive objective
 - theorem: $\log \frac{\pi^{(n)}(Y^{(n)} | x)}{\pi_0^{(n)}(Y^{(0)} | x)} - \log \frac{\pi_0(Y^{(n)} | x)}{\pi_0(Y^{(0)} | x)} = \beta_n^*$
 - $\operatorname{argmin}_{\pi} \mathbb{E}[(\log \frac{\pi(Y^{(n)} | x)}{\pi(Y^{(0)} | x)} - \log \frac{\pi_0(Y^{(n)} | x)}{\pi_0(Y^{(0)} | x)} - \beta_n^*)^2]$
 - problem: only controls *ratio*
- BonBon alignment: use both.
 - Notice: KL vs win-rate *implicitly* controlled by n .

Goal

Align LLM to target policy *equal to its own best-of- n sampling distribution*

BonBon Alignment

- Idea: use supervised finetuning (MLE) on best-of- n samples.
 - Problem: very slow.
- Idea: use best-of- n and worst-of- n samples to define contrastive objective
 - theorem: $\log \frac{\pi^{(n)}(Y^{(n)} | x)}{\pi_0^{(n)}(Y^{(0)} | x)} - \log \frac{\pi_0(Y^{(n)} | x)}{\pi_0(Y^{(0)} | x)} = \beta_n^*$
 - $\operatorname{argmin}_{\pi} \mathbb{E}[(\log \frac{\pi(Y^{(n)} | x)}{\pi(Y^{(0)} | x)} - \log \frac{\pi_0(Y^{(n)} | x)}{\pi_0(Y^{(0)} | x)} - \beta_n^*)^2]$
 - problem: only controls *ratio*
- BonBon alignment: use both.
 - Notice: KL vs win-rate *implicitly* controlled by n .

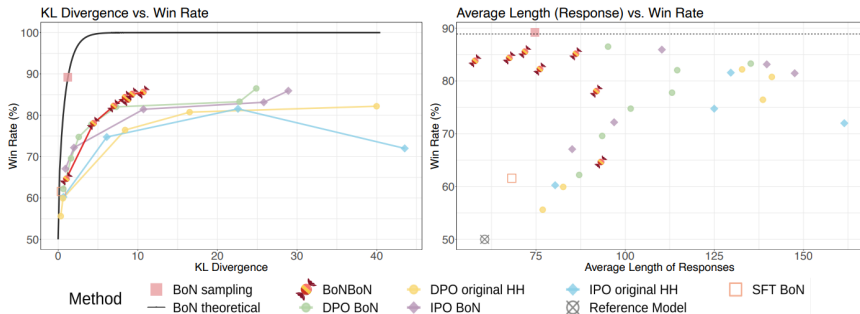
Goal

Align LLM to target policy *equal to its own best-of- n sampling distribution*

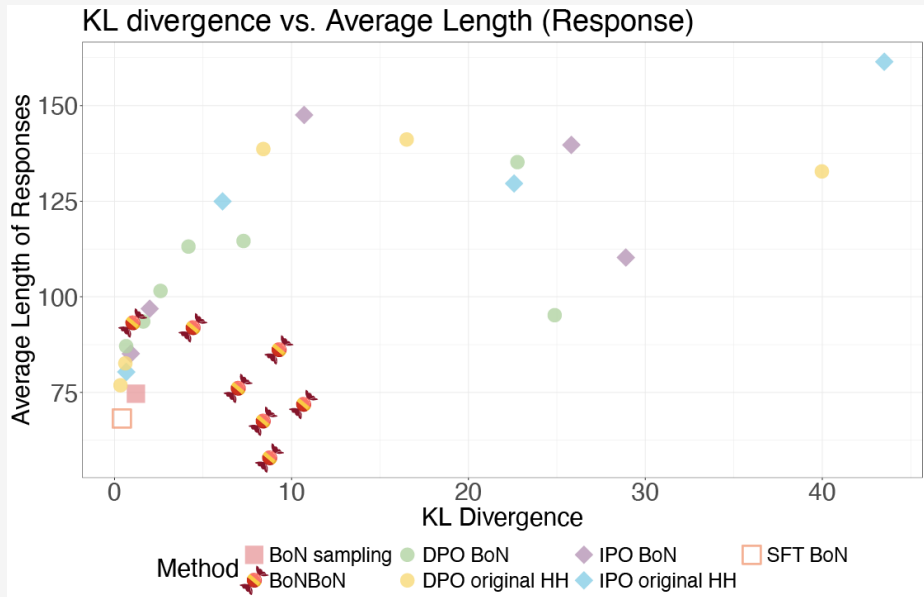
BonBon Alignment

- Idea: use supervised finetuning (MLE) on best-of- n samples.
 - Problem: very slow.
- Idea: use best-of- n and worst-of- n samples to define contrastive objective
 - theorem: $\log \frac{\pi^{(n)}(Y^{(n)} | x)}{\pi_0^{(n)}(Y^{(0)} | x)} - \log \frac{\pi_0(Y^{(n)} | x)}{\pi_0(Y^{(0)} | x)} = \beta_n^*$
 - $\operatorname{argmin}_{\pi} \mathbb{E}[(\log \frac{\pi(Y^{(n)} | x)}{\pi(Y^{(0)} | x)} - \log \frac{\pi_0(Y^{(n)} | x)}{\pi_0(Y^{(0)} | x)} - \beta_n^*)^2]$
 - problem: only controls *ratio*
- BonBon alignment: use both.
 - Notice: KL vs win-rate *implicitly* controlled by n .

Helpful and Harmless



Measured KL is Deceptive



RATE: Score Reward Models with Imperfect Rewrites of Rewrites

Question

Does a given reward model R actually reward an off-target behavior?

- Naive: measure correlation between reward and off-target behavior
Problem: could have spurious correlation
- Better: *rewrite* responses to change concept, then compare original and rewrite rewards
Problem: imperfect rewriting might change many things
- Idea: compare reward of rewrite to reward of *rewrite-of-rewrite*

Question

Does a given reward model R actually reward an off-target behavior?

- Naive: measure correlation between reward and off-target behavior
Problem: could have spurious correlation
- Better: *rewrite* responses to change concept, then compare original and rewrite rewards
Problem: imperfect rewriting might change many things
- Idea: compare reward of rewrite to reward of *rewrite-of-rewrite*

Question

Does a given reward model R actually reward an off-target behavior?

- Naive: measure correlation between reward and off-target behavior
Problem: could have spurious correlation
- Better: *rewrite* responses to change concept, then compare original and rewrite rewards
Problem: imperfect rewriting might change many things
- Idea: compare reward of rewrite to reward of *rewrite-of-rewrite*

Question

Does a given reward model R actually reward an off-target behavior?

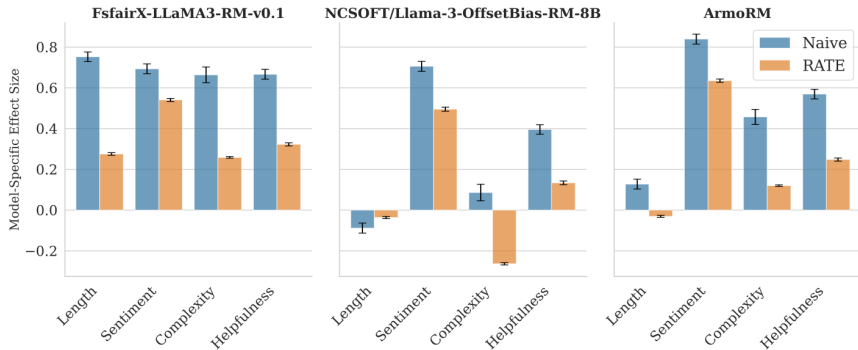
- Naive: measure correlation between reward and off-target behavior
Problem: could have spurious correlation
- Better: *rewrite* responses to change concept, then compare original and rewrite rewards
Problem: imperfect rewriting might change many things
- Idea: compare reward of rewrite to reward of *rewrite-of-rewrite*

Rewrite-of-Rewrites

Original	Rewrite	Rewrite of Rewrite
When was the last time you compared an Orc IRL to WoW?	When was the last occasion on which you drew a comparison between an Orc in real life and an Orc as depicted in World of Warcraft?	When did you last compare a real-life Orc to a World of Warcraft Orc?
W = 0, Reward: 0.14	W = 1, Reward: 0.12	W = 0, Reward: 0.16
Pros for ssd's: -Smaller form factors available - Significantly faster read- /write speeds -Very low th...	Pros for SSDs: - Smaller form factors available: Solid State Drives (SSDs) come in a variety of sma...	Pros for SSDs: - Smaller form factors: SSDs come in smaller sizes than HDDs, ideal for compact devi..
W = 0, Reward: 0.13	W = 1, Reward: 0.17	W = 0, Reward: 0.16
It wouldn't make things better; you would just end up with a hurricane full of radioactive dust and ...	Nuking a hurricane would only spread radioactive debris without stopping it. Two key points: First, ...	Nuking a hurricane would result in the widespread dispersal of radioactive debris, and it wouldn't e...
W = 1, Reward: 0.135	W = 0, Reward: 0.134	W = 1, Reward: 0.139

Table 4: Whether for a rewrite or a rewrite-of-a-rewrite, GPT-4o uses well-formatted text and a slightly formal tone. Here, W is length; samples are drawn from the ELI5 dataset, scored using ArmoRM, and truncated to 100 characters for display. The first was selected for illustrative purposes, the latter two were randomly selected from the dataset.

Naive vs RATE Estimates Across Models



BonBon Alignment

- Best-of- n shows alignment with minimal off-target drift is possible.
- RATE shows reward models have only moderate off-target bias.

Upshot

much of off-target drift appears to be a methodology bug



Lin Gui



Cristina Garbacea

BoNBoN Alignment for Large Language Models and the Sweetness of Best-of- n



David Reber



Sean Richardson



Todd Nief



Cristina Garbacea

RATE: Score Reward Models with Imperfect Rewrites of Rewrites