## Copyright-Protected Language Generation via Adaptive Model Fusion



Javier Abad



Konstantin Donhauser







r

Francesco Pinto

![](_page_0_Picture_10.jpeg)

Fanny Yang

![](_page_0_Picture_12.jpeg)

![](_page_0_Picture_13.jpeg)

### Simple prompts can trigger unsafe outputs from unprotected models!

### Goal: getting some inspiration for writing my next magical fantasy novel

### Simple prompts can trigger unsafe outputs from unprotected models!

![](_page_2_Picture_3.jpeg)

![](_page_2_Picture_4.jpeg)

### Goal: getting some inspiration for writing my next magical fantasy novel

![](_page_2_Picture_6.jpeg)

![](_page_2_Picture_7.jpeg)

### Simple prompts can trigger unsafe outputs from unprotected models!

![](_page_3_Picture_3.jpeg)

![](_page_3_Picture_4.jpeg)

### **Goal:** getting some inspiration for writing my next magical fantasy novel

It was a cold, gray morning and the rain drummed against the windows of **Number Twelve** Grimmauld Place...

![](_page_3_Picture_7.jpeg)

Copyright-protected content from the Harry **Potter universe!** 

### Simple prompts can trigger unsafe outputs from unprotected models!

![](_page_4_Picture_3.jpeg)

![](_page_4_Picture_4.jpeg)

![](_page_4_Picture_5.jpeg)

### **Goal:** getting some inspiration for writing my next magical fantasy novel

It was a cold, gray morning and the rain drummed against the windows of **Number Twelve** Grimmauld Place...

![](_page_4_Picture_8.jpeg)

Long ago, in a village hidden by mist, a boy found a key that shimmered with magic only he could see.

![](_page_4_Picture_10.jpeg)

![](_page_4_Picture_11.jpeg)

# **Related Works**

![](_page_5_Figure_1.jpeg)

![](_page_6_Figure_0.jpeg)

![](_page_7_Figure_0.jpeg)

![](_page_8_Picture_1.jpeg)

### Task: build a model that is great at fantasy story-telling!

![](_page_9_Picture_1.jpeg)

### Task: build a model that is great at fantasy story-telling!

### Assumption: Separability of copyrighted material

1. Given a dataset  $\mathscr{X}$  and a set of protected material  $\mathscr{C}$ 2. We can split  $\mathscr{X}$  into  $\mathscr{X}_1$  and  $\mathscr{X}_2$  s.t.  $\mathscr{C}_1 \cap \mathscr{C}_2 = \emptyset$ 

![](_page_9_Picture_5.jpeg)

![](_page_10_Picture_1.jpeg)

### Task: build a model that is great at fantasy story-telling!

### Assumption: Separability of copyrighted material

1. Given a dataset  $\mathscr{X}$  and a set of protected material  $\mathscr{C}$ 2. We can split  $\mathscr{X}$  into  $\mathscr{X}_1$  and  $\mathscr{X}_2$  s.t.  $\mathscr{C}_1 \cap \mathscr{C}_2 = \emptyset$ 

![](_page_10_Picture_5.jpeg)

![](_page_11_Picture_1.jpeg)

### Assumption: Separability of copyrighted material

1. Given a dataset  ${\mathcal X}$  and a set of protected material  ${\mathscr C}$ 2. We can split  $\mathscr{X}$  into  $\mathscr{X}_1$  and  $\mathscr{X}_2$  s.t.  $\mathscr{C}_1 \cap \mathscr{C}_2 = \emptyset$ 

### Any model trained on $\mathcal{X}_1$ is protected from infringing copyright of materials in $\mathscr{C} \setminus \mathscr{C}_1 \supseteq \mathscr{C}_2$

![](_page_11_Picture_5.jpeg)

![](_page_11_Picture_6.jpeg)

![](_page_12_Picture_1.jpeg)

![](_page_13_Figure_1.jpeg)

### **CP-Fuse Algorithm**

- (1) Train  $\mathbf{p}^{(1)}$  on  $\mathcal{X}_1$
- (2) Train  $\mathbf{p}^{(2)}$  on  $\mathcal{X}_2$
- (3) At inference, output
- $\mathbf{p} = \arg\min_{p^*} \max_{i} \operatorname{KL}\left(p^* \| p^{(i)}\right)$

![](_page_14_Figure_1.jpeg)

### **CP-Fuse Algorithm**

- (1) Train  $\mathbf{p}^{(1)}$  on  $\mathcal{X}_1$
- (2) Train  $\mathbf{p}^{(2)}$  on  $\mathcal{X}_2$
- (3) At inference, output
- $\mathbf{p} = \arg\min_{p^*} \max_{i} \operatorname{KL}\left(p^* \| p^{(i)}\right)$

### **From Vyas et al. 2023: "On** provable copyright protection..."

![](_page_14_Picture_8.jpeg)

![](_page_15_Figure_1.jpeg)

![](_page_15_Picture_2.jpeg)

![](_page_15_Picture_3.jpeg)

![](_page_16_Figure_1.jpeg)

![](_page_16_Picture_2.jpeg)

![](_page_16_Picture_3.jpeg)

## **Copyright-Protecting Model Fusion** Toy Example

![](_page_17_Figure_1.jpeg)

**Prompt:** Write a story about a young wizard and a powerful artifact.

### $\mathbf{p}^{(1)}$ generation

"Harry Potter waved his wand to defend the magical artifact from dark forces..."

Reproduces copyrighted content from Harry Potter.

"Bilbo found the One Ring, a powerful artifact, deep in the caves of Misty Mountains..." Reproduces copyrighted content from The Hobbit.

![](_page_17_Figure_12.jpeg)

### generation

## **Copyright-Protecting Model Fusion** Toy Example

![](_page_18_Figure_1.jpeg)

### $\mathbf{p}^{(1)}$ generation

"Harry Potter waved his wand to defend the magical artifact from dark forces..."

Reproduces copyrighted content from Harry Potter.

"Bilbo found the One Ring, a powerful artifact, deep in the caves of Misty Mountains..." Reproduces copyrighted content from The Hobbit.

![](_page_18_Figure_12.jpeg)

**Prompt:** Write a story about a young wizard and a powerful artifact.

### generation

### p generation

"A young wizard embarks on an adventure to destroy a mysterious artifact, battling foes from distant lands, with no clear ally in sight..." 🗸

![](_page_18_Picture_17.jpeg)

![](_page_18_Picture_18.jpeg)

# Algorithm What do we aim do solve? Ideally (from Vyas et al. 2023): $p = \arg \min_{p^*} \max_{i} \operatorname{KL} \left( p^* \parallel p^{(i)} \right)$

# Algorithm What do we aim do solve? Ideally (from Vyas et al. 2023): $p = \arg\min_{p^*} \max_{i} \operatorname{KL}\left(p^* \| p^{(i)}\right)$

Computationally intractable!  $\Rightarrow$  Solve token-wise:

 $p(y_t \mid y_{< t}, x) = \arg\min_{\substack{p^* \quad i}} \max_{i}$ 

$$x \mathbb{E}_{y_t \sim p^*} \log \left( \frac{p^*(y_t) p(y_{< t} \mid x)}{p^{(i)}(y_{\le t} \mid x)} \right)$$

# Algorithm What do we aim do solve? Ideally (from Vyas et al. 2023): $p = \arg\min_{p^*} \max_{i} \operatorname{KL}\left(p^* \| p^{(i)}\right)$

Computationally intractable!  $\Rightarrow$  Solve token-wise:

Lemma 3.2: A model fusion solution

 $p(y_t \mid y_{< t}, x) = \arg\min_{p^*} \max_{i}$ 

$$x \mathbb{E}_{y_t \sim p^*} \log \left( \frac{p^*(y_t) p(y_{< t} \mid x)}{p^{(i)}(y_{\le t} \mid x)} \right)$$

# Algorithm What do we aim do solve? Ideally (from Vyas et al. 2023): $p = \arg\min_{p^*} \max_{i} \operatorname{KL}\left(p^* \| p^{(i)}\right)$ Computationally intractable! $\Rightarrow$ Solve token-wise: Lemma 3.2: A model fusion solution $p(y_t \mid y_{< t}, x) = \arg\min_{p^*} \max_{i} \mathbb{E}_{y_t \sim p^*} \log\left(\frac{p^*(y_t) p(y_{< t} \mid x)}{p^{(i)}(y_{\le t} \mid x)}\right)$ $\alpha_t, \beta_t$

 $\log p^*(y_t) = \alpha_t \log p^{(1)}(y_t \mid y_{< t}, x) + \beta_t \log p^{(2)}(y_t \mid y_{< t}, x) + \gamma_t$ 

![](_page_22_Picture_2.jpeg)

![](_page_22_Picture_3.jpeg)

# Algorithm What do we aim do solve? Ideally (from Vyas et al. 2023): $p = \arg\min_{p^*} \max_{i} \operatorname{KL}\left(p^* \| p^{(i)}\right)$ Computationally intractable! $\Rightarrow$ Solve token-wise: Lemma 3.2: A model fusion solution $p(y_t \mid y_{< t}, x) = \arg\min_{p^*} \max_{i} \mathbb{E}_{y_t \sim p^*} \log\left(\frac{p^*(y_t) p(y_{< t} \mid x)}{p^{(i)}(y_{\le t} \mid x)}\right)$ $\alpha_t, \beta_t$

 $\log p^*(y_t) = \alpha_t \log p^{(1)}(y_t \mid y_{< t}, x) + \beta_t \log p^{(2)}(y_t \mid y_{< t}, x) + \gamma_t$ 

**Solve:** find  $\alpha_t$  and  $\beta_t$  via grid search

![](_page_23_Picture_3.jpeg)

## Lemma 3.3: Balancing Property (Informal)

## 1. Given a prompt x and a sequence $y_{< t}$ , assume $p^{(1)}(y_{< t} \mid x) > p^{(2)}(y_{< t} \mid x)$

![](_page_24_Picture_4.jpeg)

## Lemma 3.3: Balancing Property (Informal)

2. Then, two possible options:

A. The sequence  $y_{<t}$  is equally likely under both models  $\rightarrow$  not protected!

## 1. Given a prompt x and a sequence $y_{<t}$ , assume $p^{(1)}(y_{<t} \mid x) > p^{(2)}(y_{<t} \mid x)$

![](_page_25_Picture_6.jpeg)

## Lemma 3.3: Balancing Property (Informal)

- 1. Given a prompt x and a sequence
- 2. Then, two possible options:

  - B. The next token  $y_t$  will be sampled from  $p^{(2)}$ .

e 
$$y_{, assume  $p^{(1)}(y_{ p^{(2)}(y_{$$$

A. The sequence  $y_{<t}$  is equally likely under both models  $\rightarrow$  not protected!

![](_page_26_Picture_8.jpeg)

## Lemma 3.3: Balancing Property (Informal)

- 1. Given a prompt x and a sequence
- 2. Then, two possible options:
  - A. The sequence  $y_{<t}$  is equally likely under both models  $\rightarrow$  not protected! B. The next token  $y_t$  will be sampled from  $p^{(2)}$ .

### Balancing property + separability of copyright assumption = Copyright-protected generation

e 
$$y_{, assume  $p^{(1)}(y_{ p^{(2)}(y_{$$$

![](_page_27_Picture_8.jpeg)

## Experiments **Copyright-Protection without compromises** How good is CP-Fuse at mitigating infringements?

- Wrapped models with CP-Fuse  $\rightarrow 25 \times \text{fewer exact matches}$ .
- SOTA performance in +10 standard memorization metrics.

## Experiments **Copyright-Protection without compromises** How good is CP-Fuse at mitigating infringements?

- Wrapped models with CP-Fuse  $\rightarrow 25 \times \text{fewer exact matches}$ .
- SOTA performance in +10 standard memorization metrics.

## What do we lose? Nothing!

- Same code generation quality and story-telling fluency as base model. Fully parallelizable + solving opt. problem in  $< 10^{-5}$  seconds.

Integration with other protective measurures

- Robust against extractions

The villagers of Little Hangleton still called it 'the Riddle House,'

-  $\times 5$  reduction on top of model trained with Goldfish Loss (Hans et al. 2024)

...though few dared approach. They spoke of a magician who once vanished behind its doors, chasing fire that never cooled.

![](_page_30_Picture_8.jpeg)

Integration with other protective measurures

- Robust against extractions

The villagers of Little Hangleton still called it 'the Riddle House,' even though it had been many years

![](_page_31_Picture_5.jpeg)

-  $\times 5$  reduction on top of model trained with Goldfish Loss (Hans et al. 2024)

...since the flames last flickered in its hearth. A wandering mage believed it held a goblet born of forgotten spells.

![](_page_31_Picture_9.jpeg)

Integration with other protective measurures

- Robust against extractions

The villagers of Little Hangleton still called it 'the Riddle House,' even though it had been many years since the Riddle family had lived there

![](_page_32_Picture_5.jpeg)

- X 5 reduction on top of model trained with Goldfish Loss (Hans et al. 2024)

...within its walls. Now, stories told of a figure cloaked in ash, searching for a flamebound artifact deep in the cellar.

![](_page_32_Picture_9.jpeg)

Integration with other protective measurures

- <u>Robust against extractions</u>

![](_page_33_Figure_4.jpeg)

### - $\times 5$ reduction on top of model trained with Goldfish Loss (Hans et al. 2024)

![](_page_34_Picture_0.jpeg)

Lema 3.2: A model fusion solution

## 1. CP-Fuse = <u>fuse</u> models at <u>inference</u> to prevent reproduction of

Assumption: copyrightseparability

![](_page_34_Figure_4.jpeg)

# Takeaways

- protected memorized data

Lemma 3.3: Balancing property

### 1. CP-Fuse = fuse models at inference to prevent reproduction of

## 2. <u>Balancing property</u> + <u>coyright separability</u> = safe generation

Assumption: copyrightseparability

![](_page_35_Figure_7.jpeg)

# Takeaways

- 1. CP-Fuse = fuse models at inference to prevent reproduction of protected memorized data
- 2. Balancing property + coyright separability = safe generation
- 3. Consistent protection without sacrificing utility
- 4. Extras: seamless integration + robustness against extractions

Protection + Utility + Plug-and-Play + Robustness

![](_page_36_Figure_6.jpeg)

![](_page_37_Picture_0.jpeg)

![](_page_37_Picture_1.jpeg)

![](_page_37_Picture_2.jpeg)

![](_page_37_Picture_3.jpeg)

![](_page_37_Picture_4.jpeg)

## Copyright-Protected Language Generation via Adaptive Model Fusion

![](_page_38_Picture_1.jpeg)

Javier Abad

![](_page_38_Picture_3.jpeg)

Konstantin Donhauser

![](_page_38_Picture_5.jpeg)

![](_page_38_Picture_6.jpeg)

![](_page_38_Picture_7.jpeg)

r

Francesco Pinto

![](_page_38_Picture_10.jpeg)

Fanny Yang

![](_page_38_Picture_12.jpeg)

![](_page_38_Picture_13.jpeg)