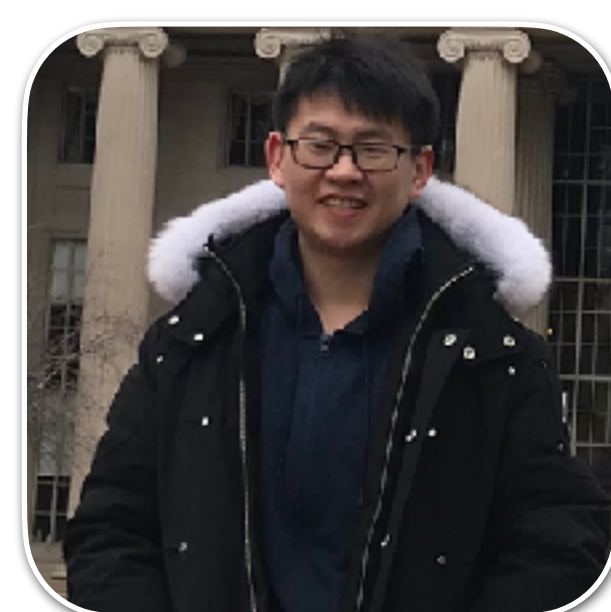# Theoretically Grounded Framework for LLM Watermarking: A Distribution-Adaptive Approach

**Haiyun He**

Postdoc @ Center for Applied Mathematics, Cornell University

**ITA 2025**

Yepeng Liu
Univ. of Florida

Prof. Ziqiao Wang
Tongji Univ.

Prof. Yongyi Mao
Univ. of Ottawa

Prof. Yuheng Bu
Univ. of Florid
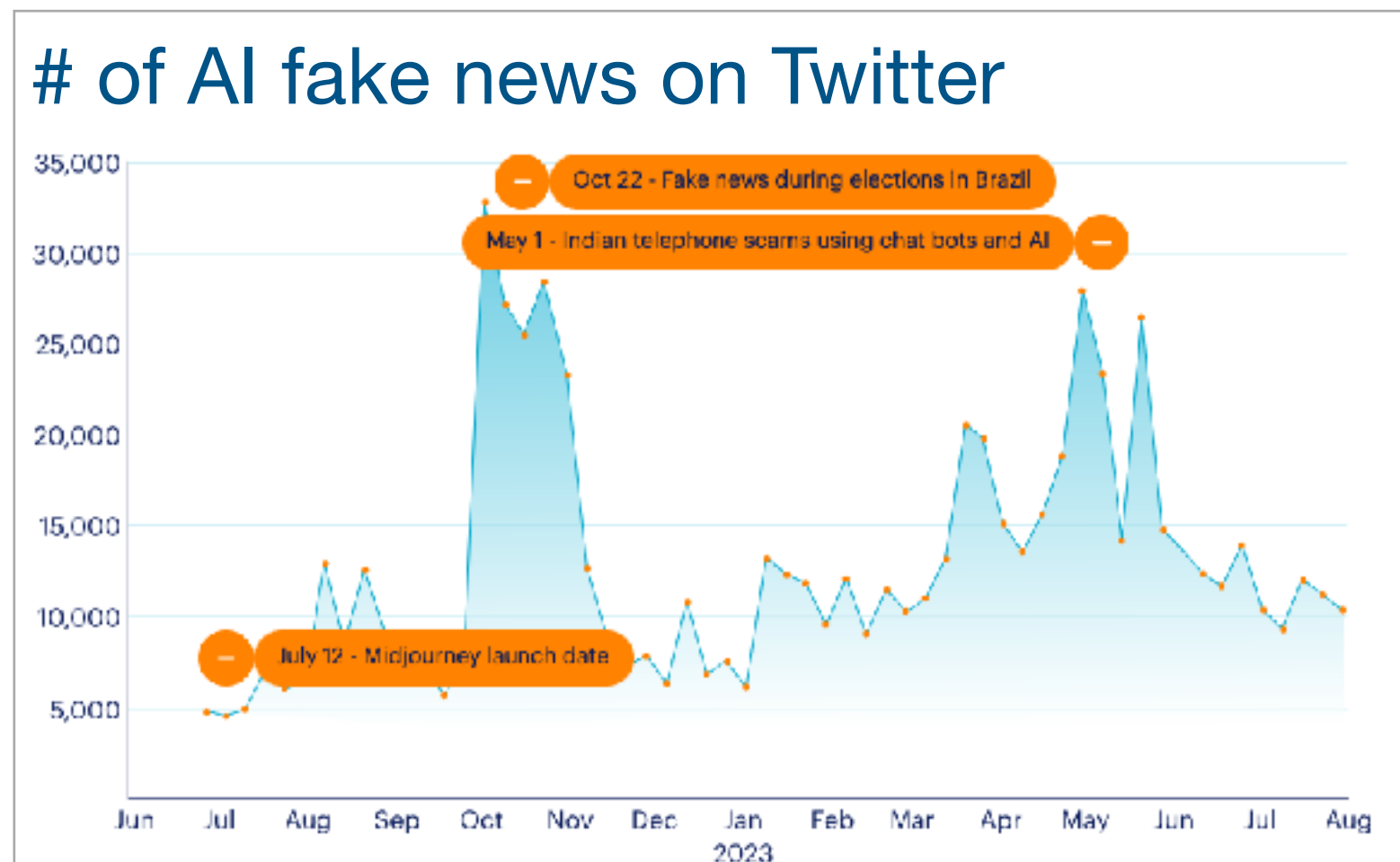
# Challenges in AI Safety

Misuse of AI-generated content

# Challenges in AI Safety

Misuse of AI-generated content



Fake news

# Challenges in AI Safety

Misuse of AI-generated content



# of AI fake news on Twitter



AI scams

# Challenges in AI Safety

## Misuse of AI-generated content

# of AI fake news on Twitter



Plagiarism

# Challenges in AI Safety

**Misuse of AI-generated content**

**Data Pollution**

# of AI fake news on Twitter



Plagiarism

# Challenges in AI Safety

## Misuse of AI-generated content

# of AI fake news on Twitter



Plagiarism

## Data Pollution

Tons of AI-generate data over the internet

# Challenges in AI Safety



**Misuse of AI-generated content**

# of AI fake news on Twitter

Plagiarism

**Data Pollution**

Tons of AI-generate data over the internet

LIVE
BREAKING NEWS
GOOGLE BANS CAT PICS ON THE INTERNET!
6:02pm APOCALYPSE MEOW!

# Challenges in AI Safety

## Misuse of AI-generated content
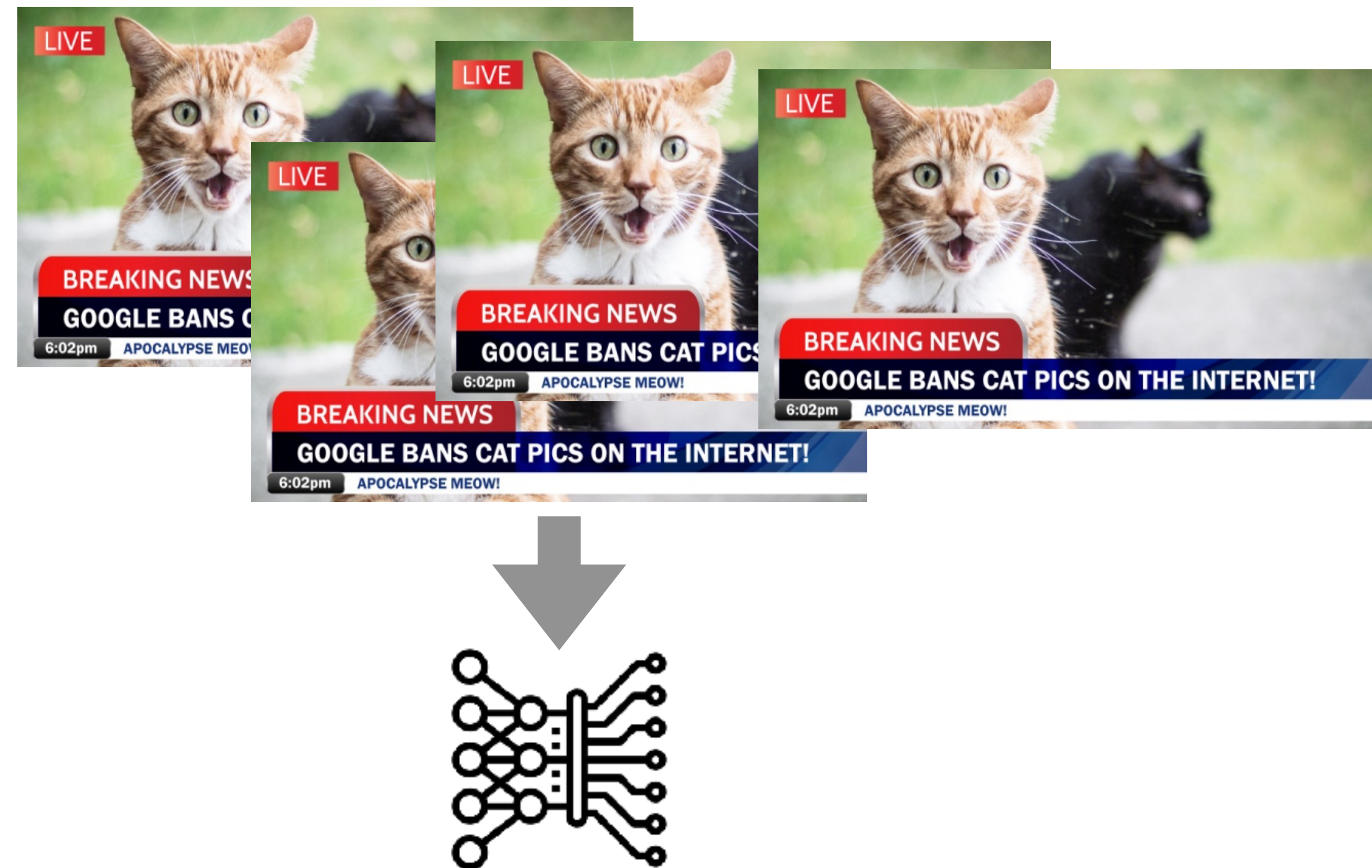
# of AI fake news on Twitter



Plagiarism

## Data Pollution

Tons of AI-generate data over the internet



???
**No cats anymore?**

# Challenges in AI Safety

## Misuse of AI-generated content

# of AI fake news on Twitter



Plagiarism

## Data Pollution

Tons of AI-generate data over the internet



???
No cats anymore?

collapse…

# Challenges in AI Safety

**Misuse of AI-generated content**

**Data Pollution**

# of AI fake news on Twitter

Tons of AI-generate data over the internet

LIVE LIVE LIVE

N THE INTERNET!

We must distinguish AI-generated data from authentic, naturally occurring data!

Plagiarism

???
No cats anymore?

collapse…

# Identify AI-generated Text

**Possible solutions?**

# Identify AI-generated Text
## Possible solutions?

- By observation:

# Identify AI-generated Text
## Possible solutions?

*"Here's the revised version of your…"*, *"Best regards,[Your Name]"*    :-D

# Identify AI-generated Text
## Possible solutions?

- Metadata  <—easy to remove

**Metadata**

**File name:** Dataset
**Author:** GPT
**Location:** Ithaca
**Created:** Jan 01, 2025

# Identify AI-generated Text
## Possible solutions?

- Giant database to store all AI-generated content <—storage? privacy?

# Identify AI-generated Text
**Possible solutions?**

- Discriminator models:  GPTZero  DetectGPT Copyleaks  pangramlabs ...

# Identify AI-generated Text
## Possible solutions?

<—high prob of falsely alarming human-written text

# Identify AI-generated Text
## Possible solutions?

- **Watermarking: inserting a signal into LLM predicted tokens**

# Identify AI-generated Text

**Possible solutions?**



- **Watermarking: inserting a signal into LLM predicted tokens**

# Identify AI-generated Text

## Possible solutions?



- **Watermarking: inserting a signal into LLM**

# Identify AI-generated Text
## Possible solutions?



- **Watermarking: inserting a signal into LLM predicted tokens**

# Identify AI-generated Text
## Possible solutions?

- **Watermarking: inserting a signal into LLM predicted tokens**

# A Framework for LLM Watermark Generation

# A Framework for LLM Watermark Generation

# A Framework for LLM Watermark Generation

Distribution of $x_t$ : $Q_{X_t|X_1^{t-1}}$

*Wolfgang Mozart was an influential*

LLM

composer
musician
Vienna
and

# A Framework for LLM Watermark Generation

Distribution of $x_t$ : $Q_{X_t|X_1^{t-1}}$

*Wolfgang Mozart was an influential*

**LLM**

composer
musician
Vienna
and

sample

*composer*

# A Framework for LLM Watermark Generation

Distribution of $x_t : Q_{X_t|X_1^{t-1}}$

| | |
|---|---|
| *Wolfgang Mozart was an influential* | → |

LLM

→

composer
musician
Vienna
and

sample → *composer*

append

# A Framework for LLM Watermark Generation

Distribution of $x_t$ : $Q_{X_t|X_1^{t-1}}$

Wolfgang Mozart
was an influential

LLM

composer
musician
Vienna
and

sample

composer

Watermarking
Scheme

# A Framework for LLM Watermark Generation

Distribution of $x_t$ : $Q_{X_t|X_1^{t-1}}$

Wolfgang Mozart was an influential

LLM

composer
musician
Vienna
and

sample

composer

Watermarking Scheme

Insert a signal $\zeta_t$

# A Framework for LLM Watermark Generation

*Wolfgang Mozart was an influential*

**LLM**

Distribution of $x_t : Q_{X_t | X_1^{t-1}}$

composer
musician
Vienna
and

sample → *composer*

**Watermarking Scheme**

Insert a signal $\zeta_t$

Altered distribution of $x_t : P_{X_t | X_1^{t-1}, \zeta_t}$

composer
musician
Vienna
and

# A Framework for LLM Watermark Generation

*Wolfgang Mozart was an influential*

**LLM**

Distribution of $x_t$ : $Q_{X_t|X_1^{t-1}}$

composer
musician
Vienna
and

sample → *composer*

**Watermarking Scheme**

Insert a signal $\zeta_t$

Altered distribution of $x_t$ : $P_{X_t|X_1^{t-1},\zeta_t}$

composer
musician
Vienna
and

sample → *musician*

# A Framework for LLM Watermark Generation

*Wolfgang Mozart was an influential*

**LLM**

Distribution of $x_t$ : $Q_{X_t|X_1^{t-1}}$

composer
musician
Vienna
and

sample → *composer*

Watermarking Scheme

Insert a signal $\zeta_t$

Altered distribution of $x_t$ : $P_{X_t|X_1^{t-1}, \zeta_t}$

composer
musician
Vienna
and

sample → *musician*

**Still a normal sentence. Imperceptible!**

# A Framework for LLM Watermark Generation

*Wolfgang Mozart was an influential*



LLM

Distribution of $x_t$ : $Q_{X_t|X_1^{t-1}}$

- composer
- musician
- Vienna
- and

sample → *composer*

Altered distribution of $x_t$ : $P_{X_t|X_1^{t-1},\zeta_t}$

- poser
- sician
- enna
- and

sample → *musician*

Like invisible Ink (Steganography)

**Still a normal sentence. Imperceptible!**

# A Framework for LLM Watermark Generation

*Wolfgang Mozart was an influential*

**LLM**

Distribution of $x_t$ : $Q_{X_t|X_1^{t-1}}$

composer
musician
Vienna
and

sample

*composer*

Watermarking Scheme

Insert a signal $\zeta_t$

Altered distribution of $x_t$ : $P_{X_t|X_1^{t-1},\zeta_t}$

composer
musician
Vienna
and

sample

*musician*

**Still a normal sentence. Imperceptible!**

# A Framework for LLM Watermark Generation

*Wolfgang Mozart was an influential*

**LLM**

Distribution of $x_t$ : $Q_{X_t|X_1^{t-1}}$

composer
musician
Vienna
and

sample

*composer*

**Watermarking Scheme**

Insert a signal $\zeta_t$

**auxiliary random variable**

Altered distribution of $x_t$ : $P_{X_t|X_1^{t-1}, \zeta_t}$

composer
musician
Vienna
and

sample

*musician*

**Still a normal sentence. Imperceptible!**

# Hypothesis Testing for LLM Watermark Detection



$$x_1^T$$

| Wolfgang Mozart was an influential | → | LLM | → | 🔑 Watermarking Scheme | → Insert signals $\zeta_1, \zeta_2, \ldots, \zeta_T$ | ⋯ → | Wolfgang Mozart was an influential musician of the Classical period. |

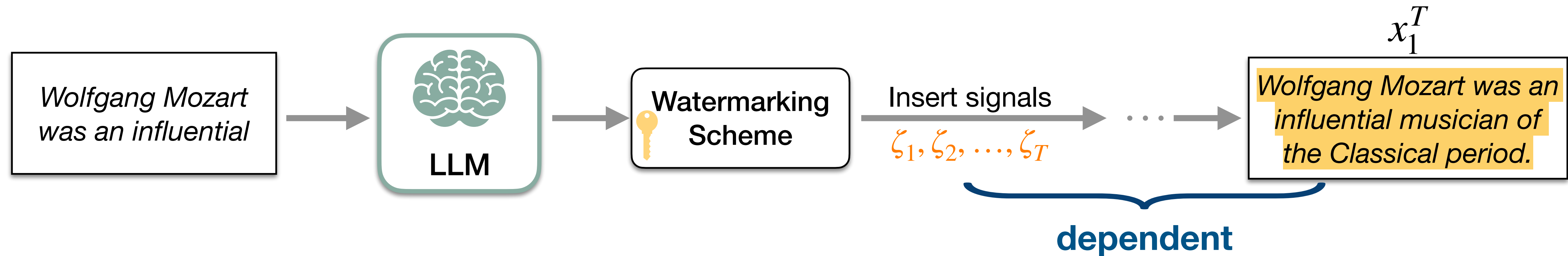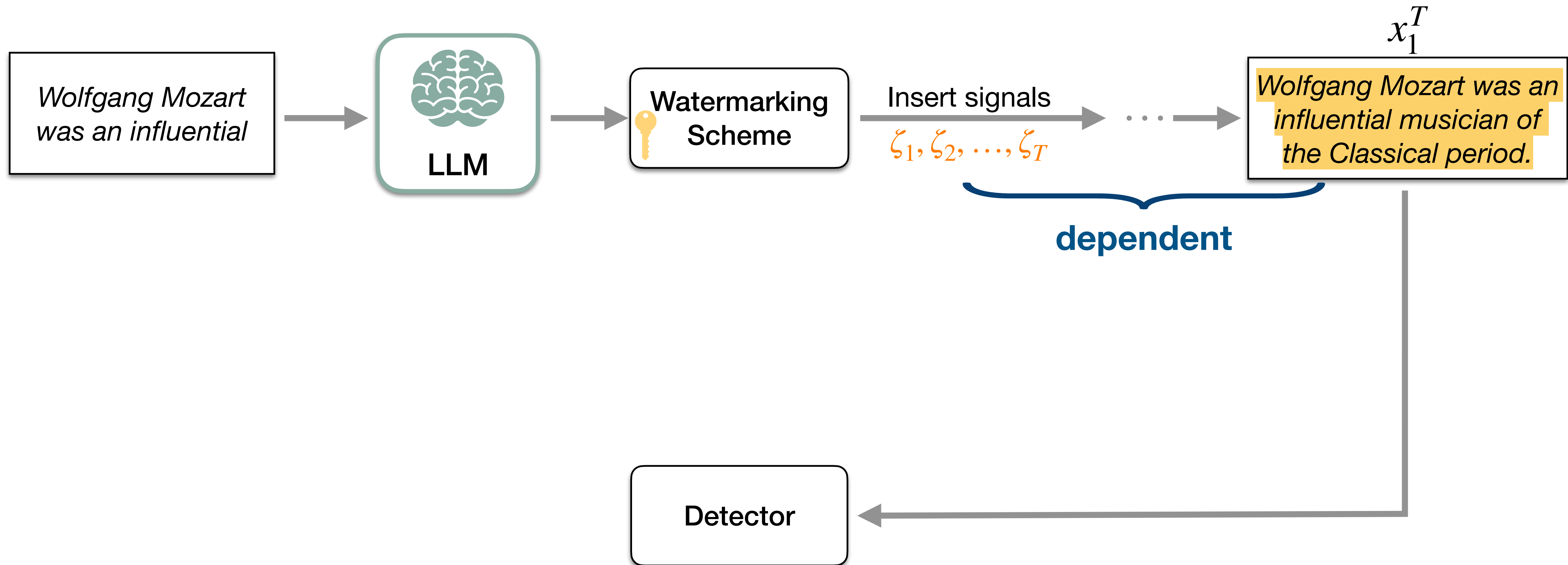# Hypothesis Testing for LLM Watermark Detection

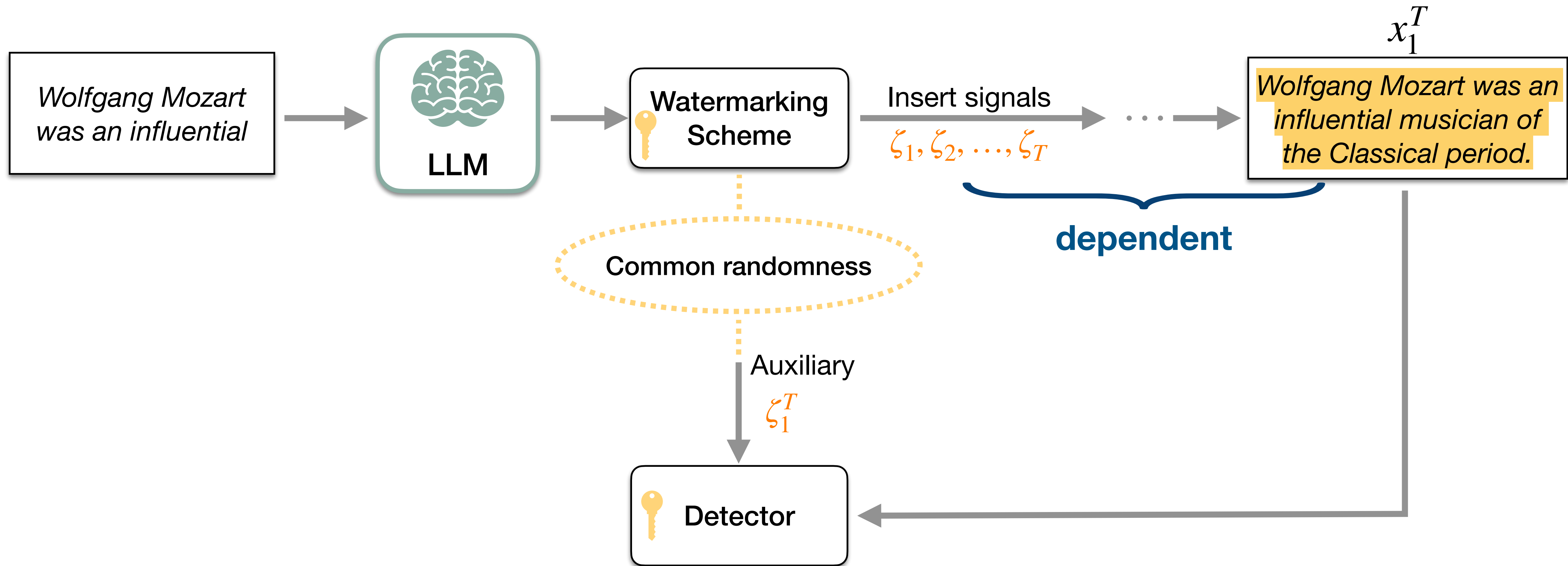# Hypothesis Testing for LLM Watermark Detection

# Hypothesis Testing for LLM Watermark Detection

# Hypothesis Testing for LLM Watermark Detection



$x_1^T$

Wolfgang Mozart was an influential

LLM

Watermarking Scheme

Insert signals
$\zeta_1, \zeta_2, \ldots, \zeta_T$

Wolfgang Mozart was an influential musician of the Classical period.

dependent

Common randomness

Auxiliary
$\zeta_1^T$

Detector

# Hypothesis Testing for LLM Watermark Detection

$x_1^T$

Wolfgang Mozart was an influential → LLM → 🔑 Watermarking Scheme → Insert signals $\zeta_1, \zeta_2, \ldots, \zeta_T$ → ···· → *Wolfgang Mozart was an influential musician of the Classical period.*

**dependent**

Common randomness

Auxiliary $\zeta_1^T$

**LLM generated** ← $x_1^T$ and $\zeta_1^T$ dependent ← 🔑 Detector

# Hypothesis Testing for LLM Watermark Detection



$x_1^T$

Wolfgang Mozart was an influential → LLM → 🔑 Watermarking Scheme → Insert signals $\zeta_1, \zeta_2, \ldots, \zeta_T$ → · · · · → Wolfgang Mozart was an influential musician of the Classical period.

**dependent**

Common randomness

Auxiliary $\zeta_1^T$

**LLM generated** ← $x_1^T$ and $\zeta_1^T$ dependent

🔑 Detector

**Human written** ← $x_1^T$ and $\zeta_1^T$ independent

# Hypothesis Testing for LLM Watermark Detection



Wolfgang Mozart was an influential

LLM

Watermarking Scheme

Insert signals $\zeta_1, \zeta_2, \ldots, \zeta_T$

$\ldots$

$x_1^T$

Wolfgang Mozart was an influential musician of the Classical period.

dependent

Common randomness

**Watermark Detection $\implies$ Hypothesis Testing:**

$$H_0 : X_1^T \text{ is human written, i.e., } (X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$$

$$H_1 : X_1^T \text{ is LLM generated, i.e., } (X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$$

# LLM Watermark Detection Errors

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

# LLM Watermark Detection Errors

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

# LLM Watermark Detection Errors

Watermark Detection $\implies$ Hypothesis Testing:   Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

# LLM Watermark Detection Errors

Watermark Detection $\implies$ Hypothesis Testing:  <span style="color:blue">Human/unwatermarked LLM</span>

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

<span style="color:red">Watermarking scheme</span>

**Performance metric:**
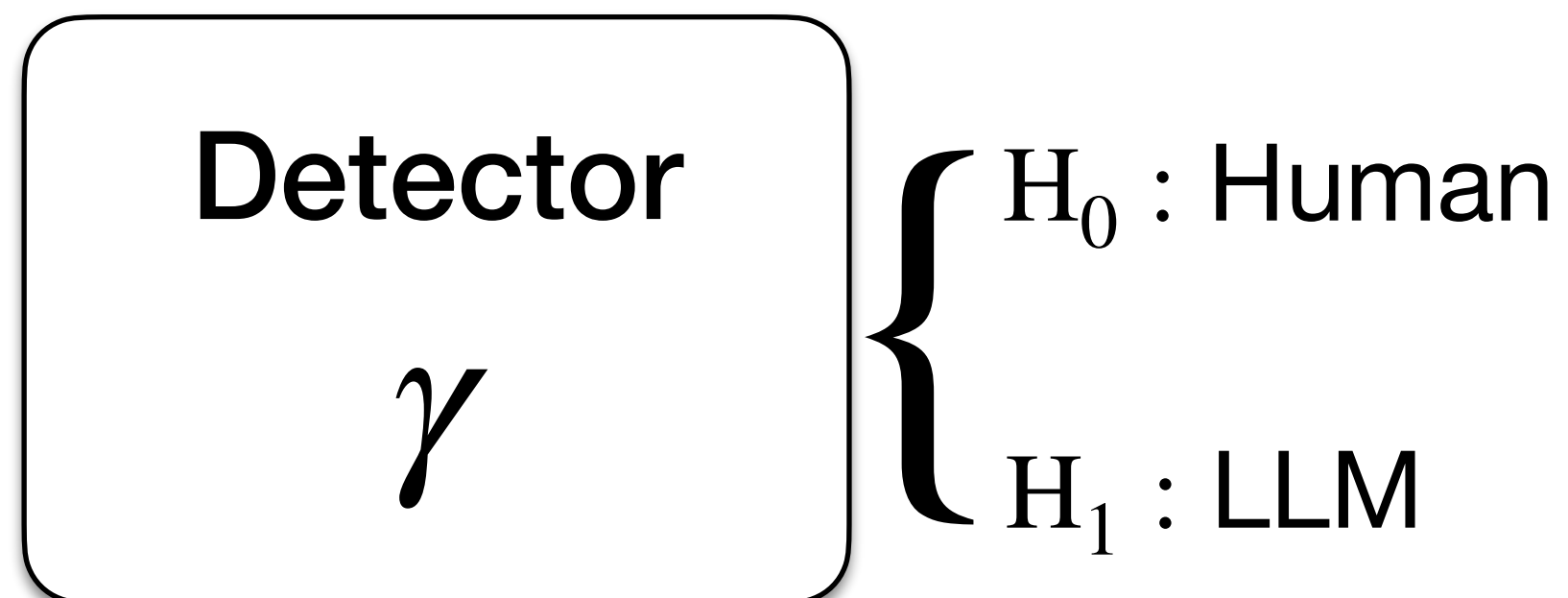
# LLM Watermark Detection Errors

Watermark Detection $\implies$ Hypothesis Testing:  Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

**<u>Performance metric:</u>**

$$\boxed{\begin{array}{c} \text{Detector} \\ \gamma \end{array}} \quad \begin{cases} H_0 : \text{Human} \\ \\ H_1 : \text{LLM} \end{cases}$$

# LLM Watermark Detection Errors

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

Watermarking scheme

**Performance metric:**

**Reality**

| | $\mathrm{H}_0$ : Human | $\mathrm{H}_1$ : LLM |
|---|---|---|
| Detector $\gamma$ $\begin{cases} \mathrm{H}_0 : \text{Human} \\ \\ \mathrm{H}_1 : \text{LLM} \end{cases}$ | | |
| | | |

# LLM Watermark Detection Errors

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$$H_0 : X_1^T \text{ is human written, i.e., } (X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$$

$$H_1 : X_1^T \text{ is LLM generated, i.e., } (X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$$

Watermarking scheme

**Performance metric:**

| | **Reality** | |
|---|---|---|
| | $H_0$ : Human | $H_1$ : LLM |
| **Detector** $\gamma$ $\begin{cases} H_0 : \text{Human} \\ \\ H_1 : \text{LLM} \end{cases}$ $\quad$ $H_0$ : Human | ✅ | |
| $H_1$ : LLM | Type-I error (false alarm) $\beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T})$ | |

# LLM Watermark Detection Errors

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

Watermarking scheme

## Performance metric:

Reality

| Detector $\gamma$ $\begin{cases} H_0 : \text{Human} \\ \\ H_1 : \text{LLM} \end{cases}$ | $H_0 : \text{Human}$ | $H_1 : \text{LLM}$ |
|---|---|---|
| | ✅ | Type-II error (miss detection) $\beta_1(\gamma, P_{X_1^T, \zeta_1^T})$ |
| | Type-I error (false alarm) $\beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T})$ | ✅ |

# LLM Watermark Detection Errors

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

**Performance metric:**

**Reality**

| Detector $\gamma$ $\begin{cases} \mathrm{H}_0 : \text{Human} \\ \\ \mathrm{H}_1 : \text{LLM} \end{cases}$ | $\mathrm{H}_0$ : Human | $\mathrm{H}_1$ : LLM |
|---|---|---|
| | ✅ | Type-II error (miss detection) $\min \ \beta_1(\gamma, P_{X_1^T, \zeta_1^T})$ |
| | Type-I error (false alarm) $\beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \ \le \alpha$ | ✅ |

# LLM Watermarked Text Quality

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

# LLM Watermarked Text Quality

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$\text{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$\text{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

scheme

Other criteria for LLM watermarking?

# LLM Watermarked Text Quality

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

scheme

Other criteria for LLM watermarking?
$\implies$ **Text Quality!**

# LLM Watermarked Text Quality

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

...cheme

Other criteria for LLM watermarking?
$\implies$ **Text Quality!**

$\implies$ **Indistinguishable from unwatermarked**

# LLM Watermarked Text Quality

Watermark Detection $\implies$ Hypothesis Testing:   Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

watermarked text distribution

$$P_{X_1^T}$$

# LLM Watermarked Text Quality

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

watermarked text distribution
$$P_{X_1^T}$$

**vs**

original text distribution
$$Q_{X_1^T}$$

7/18

# LLM Watermarked Text Quality

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

watermarked text distribution $P_{X_1^T}$ **vs** original text distribution $Q_{X_1^T}$

Good text quality

# LLM Watermarked Text Quality

Watermark Detection $\Longrightarrow$ Hypothesis Testing:   Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

Watermarking scheme

watermarked text distribution
$$P_{X_1^T}$$

**vs**

original text distribution
$$Q_{X_1^T}$$

Good text quality $\Longleftrightarrow$ $\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$

# LLM Watermarked Text Quality

Watermark Detection $\implies$ Hypothesis Testing:  Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

Watermarking scheme
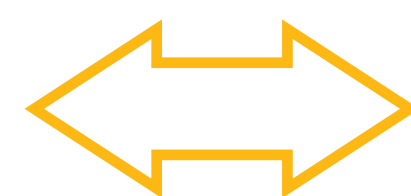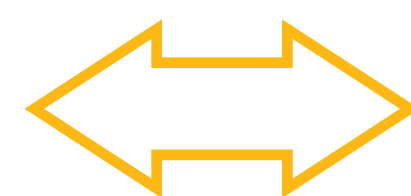
watermarked text distribution
$P_{X_1^T}$

**vs**

original text distribution
$Q_{X_1^T}$

Good text quality $\Longleftrightarrow$ $\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$

(Distortion Level)

# LLM Watermarked Text Quality

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

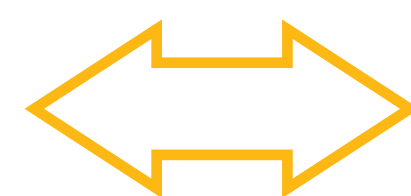Watermarking scheme

watermarked text distribution          original text distribution
$$P_{X_1^T} \qquad \textbf{vs} \qquad Q_{X_1^T}$$

Good text quality $\Longleftrightarrow$ $D(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$ (D can be any distortion metric)

(Distortion Level)

# Trade-off in Designing LLM Watermarking

Watermark Detection $\Longrightarrow$ Hypothesis Testing: Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

# Trade-off in Designing LLM Watermarking

Watermark Detection $\Longrightarrow$ Hypothesis Testing:   Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

**Trade-off:**

Type-II error——False alarm rate——Distortion Level

$$\beta_1 - - \alpha - - \epsilon$$

# Trade-off in Designing LLM Watermarking

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

# Trade-off in Designing LLM Watermarking

Watermark Detection $\implies$ Hypothesis Testing:   Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

Existing watermarking methods: **heuristic**

# Trade-off in Designing LLM Watermarking

$$\mathrm{H}_0 : X_1^T \text{ is human written, i.e., } (X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$$

$$\mathrm{H}_1 : X_1^T \text{ is LLM generated, i.e., } (X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$$

Watermarking scheme

Existing watermarking methods: **heuristic**

**Example**
[KGW-1, 2023]

No watermark
Extremely efficient on average term
lengths and word frequencies on
synthetic, microamount text (as little
as 25 words)
Very small and low-resource key/hash
(e.g., 140 bits per key is sufficient
for 99.999999999% of the Synthetic
Internet

With watermark
- minimal marginal probability for a
detection attempt.
- Good speech frequency and energy
rate reduction.
- messages indiscernible to humans.
- easy for humans to verify.

# Trade-off in Designing LLM Watermarking

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

Existing watermarking methods: **heuristic**

**Example**
[KGW-1, 2023]

No watermark
Extremely efficient on average term
lengths and word frequencies on
synthetic, microamount text (as little
as 25 words)
Very small and low-resource key/hash
(e.g., 140 bits per key is sufficient
for 99.999999999% of the Synthetic
Internet

With watermark
- minimal marginal probability for a
detection attempt.
- Good speech frequency and energy
rate reduction.
- messages indiscernible to humans.
- easy for humans to verify.

a word←randomly assign green/red color

# Trade-off in Designing LLM Watermarking

Watermark Detection $\Longrightarrow$ Hypothesis Testing:   Human/unwatermarked LLM

$$\mathrm{H}_0 : X_1^T \text{ is human written, i.e., } (X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$$

$$\mathrm{H}_1 : X_1^T \text{ is LLM generated, i.e., } (X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$$

Watermarking scheme

Existing watermarking methods: **heuristic**

**Example**
[KGW-1, 2023]

No watermark

Extremely efficient on average term
lengths and word frequencies on
synthetic, microamount text (as little
as 25 words)
Very small and low-resource key/hash
(e.g., 140 bits per key is sufficient
for 99.9999999% of the Synthetic
Internet

With watermark
- minimal marginal probability for a
detection attempt.
- Good speech frequency and energy
rate reduction.
- messages indiscernible to humans.
- easy for humans to verify.

a word←randomly assign green/red color

Green word: <u>increase</u> sampling probability

# **Trade-off in Designing LLM Watermarking**

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$\text{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\text{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

Existing watermarking methods: **heuristic**

**Example**
[KGW-1, 2023]

No watermark
Extremely efficient on average term
lengths and word frequencies on
synthetic, microamount text (as little
as 25 words)
Very small and low-resource key/hash
(e.g., 140 bits per key is sufficient
for 99.999999999% of the Synthetic
Internet

With watermark
- minimal marginal probability for a
detection attempt.
- Good speech frequency and energy
rate reduction.
- messages indiscernible to humans.
- easy for humans to verify.

○ High miss detection when requiring low false alarm

# Trade-off in Designing LLM Watermarking

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

Existing watermarking methods: **heuristic**

No watermark
Extremely efficient on average term
lengths and word frequencies on
synthetic, microamount text (as little
as 25 words)
Very small and low-resource key/hash
(e.g., 140 bits per key is sufficient
for 99.999999999% of the Synthetic
Internet

With watermark
- minimal marginal probability for a
detection attempt.
- Good speech frequency and energy
rate reduction.
- messages indiscernible to humans.
- easy for humans to verify.

High miss detection when requiring low false alarm
Not distortion-free

# Optimize LLM Watermark Generation and Detection

Watermark Detection $\Longrightarrow$ Hypothesis Testing: Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

# Optimize LLM Watermark Generation and Detection

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$\text{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\text{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

**Find the best watermarking scheme & detector:**

# Optimize LLM Watermark Generation and Detection

Watermark Detection $\implies$ Hypothesis Testing:   Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

**Find the best watermarking scheme & detector:**

Minimize miss detection $\implies$ $\min_{\gamma,\, P_{X_1^T, \zeta_1^T}} \beta_1(\gamma,\, P_{X_1^T, \zeta_1^T})$

# Optimize LLM Watermark Generation and Detection

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

Watermarking scheme

**Find the best <u>watermarking scheme</u> & <u>detector</u>:**

$$\min_{\gamma, \, P_{X_1^T, \zeta_1^T}} \beta_1(\gamma, \, P_{X_1^T, \zeta_1^T})$$

Humans are very creative, can write arbitrary texts

# Optimize LLM Watermark Generation and Detection

Watermark Detection $\Longrightarrow$ Hypothesis Testing: Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

Watermarking scheme

**Find the best watermarking scheme & detector:**

Humans are very creative, can write arbitrary texts

$$\min_{\gamma,\ P_{X_1^T,\zeta_1^T}} \beta_1(\gamma, P_{X_1^T,\zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

# Optimize LLM Watermark Generation and Detection

Watermark Detection $\implies$ Hypothesis Testing: Human/unwatermarked LLM

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

Watermarking scheme

**Find the best <u>watermarking scheme</u> & <u>detector</u>:**

$$\min_{\gamma,\ P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma,\ P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

Ensure text quality ➡

# Optimize LLM Watermark Generation and Detection

Watermark Detection $\Longrightarrow$ Hypothesis Testing:   Human/unwatermarked LLM

$\mathrm{H}_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim Q_{X_1^T} \otimes P_{\zeta_1^T}$

$\mathrm{H}_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim P_{X_1^T, \zeta_1^T}$

Watermarking scheme

**Find the best watermarking scheme & detector:**

$$\min_{\gamma,\, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma,\, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

Ensure text quality $\quad\blacktriangleright\quad$ $\mathrm{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$

# Optimize LLM Watermark Generation and Detection

Watermark Detection $\Longrightarrow$ Hypothesis Testing: **Human/unwatermarked LLM**

$H_0 : X_1^T$ is human written, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{Q_{X_1^T}} \otimes \boxed{P_{\zeta_1^T}}$

$H_1 : X_1^T$ is LLM generated, i.e., $(X_1^T, \zeta_1^T) \sim \boxed{P_{X_1^T, \zeta_1^T}}$

**Watermarking scheme**

**Find the best watermarking scheme & detector:**

$$\min_{\gamma, \, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma, \, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

# Fundamental Limit for Type-II Error

**Optimization problem:**

$$\min_{\gamma,\ P_{X_1^T,\zeta_1^T}} \quad \beta_1(\gamma,\ P_{X_1^T,\zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \le \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \le \epsilon$$

# Fundamental Limit for Type-II Error

Watermarked text distribution: $P^*_{X_1^T} = \arg \min_{P_{X_1^T}: D(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$

✦ **Minimum Type-II error:**

$$\beta_1^*(Q_{X_1^T}, \alpha, \epsilon) = \sum_{x_1^T} \left( P^*_{X_1^T}(x_1^T) - \alpha \right)_+$$

**Optimization problem:**

$$\min_{\gamma, \, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma, \, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$D(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

# Fundamental Limit for Type-II Error

Watermarked text distribution: $P^*_{X_1^T} = \arg \min\limits_{P_{X_1^T}: D(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum\limits_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$

**Optimization problem:**

$$\min\limits_{\gamma, \, P_{X_1^T, \zeta_1^T}} \beta_1(\gamma, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup\limits_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$D(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

✦ **Minimum Type-II error:**

$$\beta_1^*(Q_{X_1^T}, \alpha, \epsilon) = \sum\limits_{x_1^T} \left( P^*_{X_1^T}(x_1^T) - \alpha \right)_+$$

**Best achievable for any watermarking methods**

Same as Huang et al. (2023, Theorem 3.2) but under different framework

# Fundamental Limit for Type-II Error

Watermarked text distribution: $P^*_{X_1^T} = \arg \min_{P_{X_1^T}: D(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$

**✦ Minimum Type-II error:**

$$\beta_1^*(Q_{X_1^T}, \alpha, \epsilon) = \sum_{x_1^T} \left( P^*_{X_1^T}(x_1^T) - \alpha \right)_+$$

## **Optimization problem:**

$$\min_{\gamma, \, P_{X_1^T, \zeta_1^T}} \beta_1(\gamma, \, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$D(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

# Fundamental Limit for Type-II Error

Watermarked text distribution: $P^*_{X_1^T} = \arg \min_{P_{X_1^T}:\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$

**✦ Minimum Type-II error:**

$$\beta_1^*(Q_{X_1^T}, \alpha, \epsilon) = \sum_{x_1^T} \left( P^*_{X_1^T}(x_1^T) - \alpha \right)_+$$

**Optimization problem:**

$$\min_{\gamma,\, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma,\, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

$\mathsf{D}_{\mathsf{TV}}$

# Fundamental Limit for Type-II Error

Watermarked text distribution: $P^*_{X_1^T} = \arg \min_{P_{X_1^T} : \mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$

**Optimization problem:**

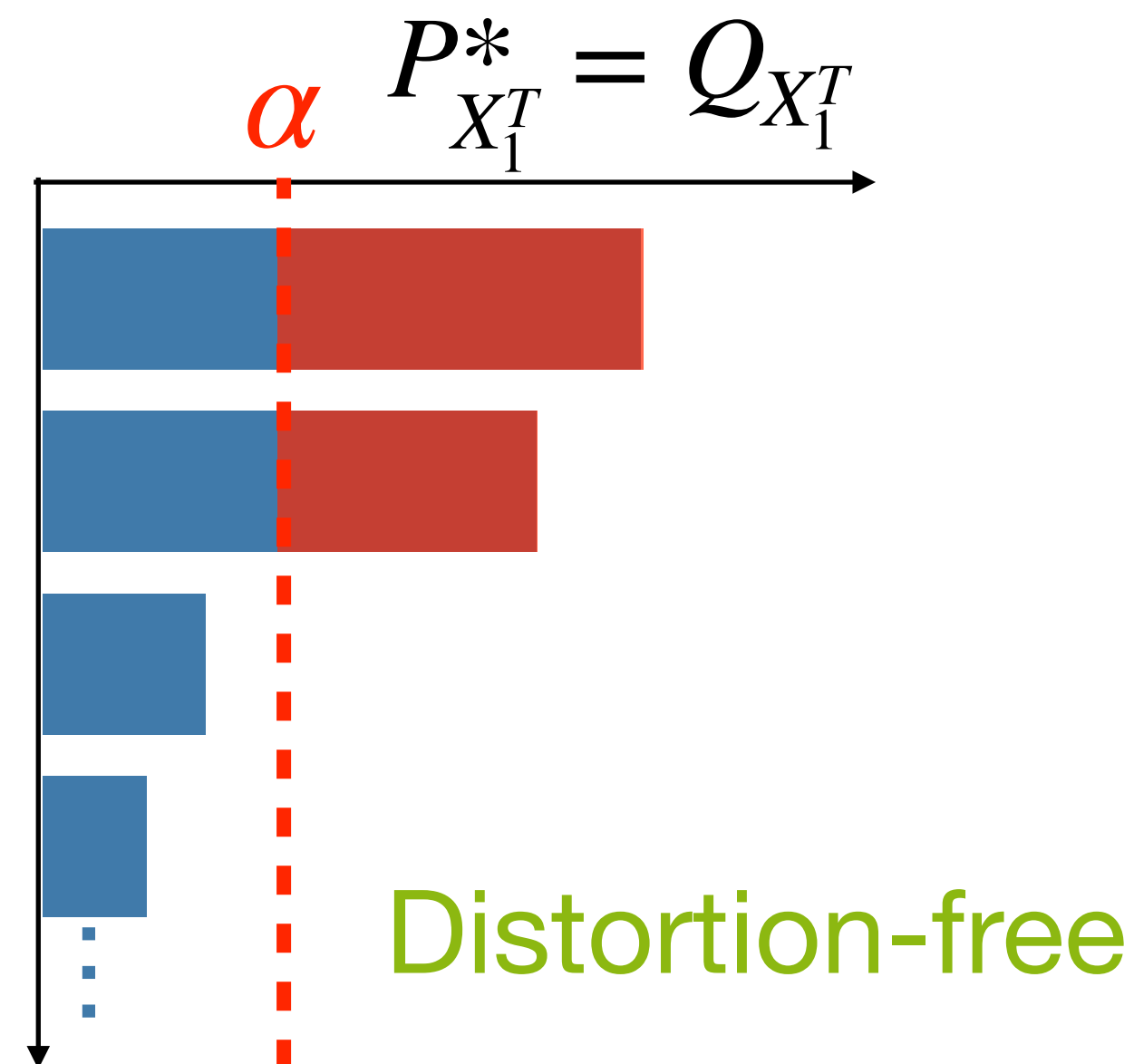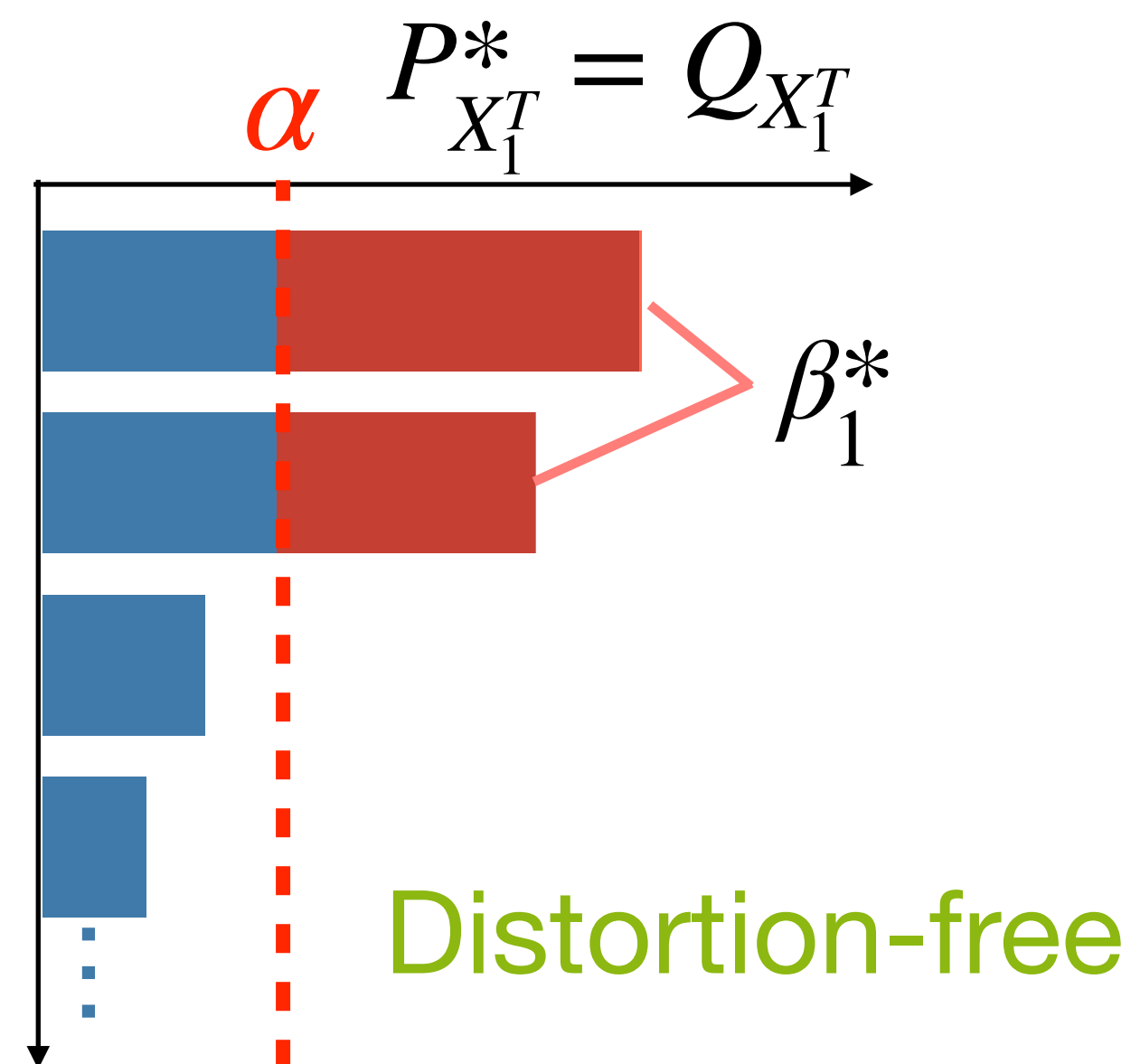$$\min_{\gamma, \, P_{X_1^T, \zeta_1^T}} \beta_1(\gamma, \, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

$\mathsf{D}_{\mathsf{TV}}$

✦ **Minimum Type-II error:**

$$\beta_1^*(Q_{X_1^T}, \alpha, \epsilon) = \sum_{x_1^T} \left( P^*_{X_1^T}(x_1^T) - \alpha \right)_+$$

$\alpha \qquad P^*_{X_1^T} = Q_{X_1^T}$

Distortion-free

# Fundamental Limit for Type-II Error

Watermarked text distribution: $P^*_{X_1^T} = \arg \min\limits_{P_{X_1^T} : D(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum\limits_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$

**Optimization problem:**

$$\min\limits_{\gamma,\, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma,\, P_{X_1^T, \zeta_1^T})$$
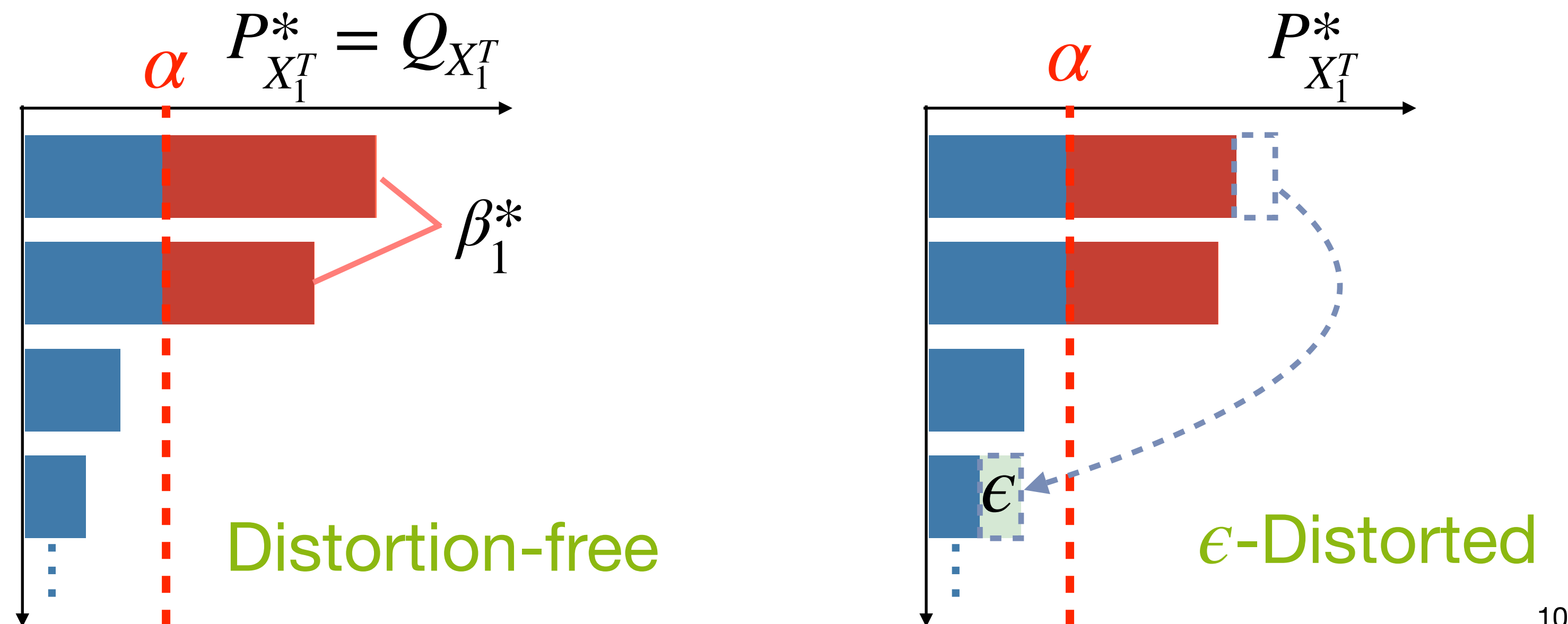
$$\text{s.t.} \quad \sup\limits_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$D(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

$D_{TV}$

✦ **Minimum Type-II error:**

$$\beta_1^*(Q_{X_1^T}, \alpha, \epsilon) = \sum\limits_{x_1^T} \left( P^*_{X_1^T}(x_1^T) - \alpha \right)_+$$



$P^*_{X_1^T} = Q_{X_1^T}$

$\alpha$

$\beta_1^*$

Distortion-free

# Fundamental Limit for Type-II Error

Watermarked text distribution: $P^*_{X_1^T} = \arg \min\limits_{P_{X_1^T}: \mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \le \epsilon} \sum\limits_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$

✦ **Minimum Type-II error:**

$$\beta_1^*(Q_{X_1^T}, \alpha, \epsilon) = \sum\limits_{x_1^T} \left(P^*_{X_1^T}(x_1^T) - \alpha\right)_+$$

## Optimization problem:

$$\min\limits_{\gamma, \, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma, \, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup\limits_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \le \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \le \epsilon$$

$\mathsf{D}_{\mathsf{TV}}$



Distortion-free

$\epsilon$-Distorted

# Fundamental Limit for Type-II Error

Watermarked text distribution: $P^*_{X^T_1} = \arg \min\limits_{P_{X^T_1}:\mathsf{D}(P_{X^T_1},Q_{X^T_1}) \leq \epsilon} \sum_{x^T_1} (P_{X^T_1}(x^T_1) - \alpha)_+$

## Optimization problem:

$$\min\limits_{\gamma,\, P_{X^T_1,\zeta^T_1}} \beta_1(\gamma,\, P_{X^T_1,\zeta^T_1})$$

s.t. $\sup\limits_{Q_{X^T_1}} \beta_0(\gamma, Q_{X^T_1}, P_{\zeta^T_1}) \leq \alpha$

$\mathsf{D}(P_{X^T_1}, Q_{X^T_1}) \leq \epsilon$

$\mathsf{D}_{\mathsf{TV}}$

✦ **Minimum Type-II error:**

$$\beta^*_1(Q_{X^T_1}, \alpha, \epsilon) = \sum_{x^T_1} \left( P^*_{X^T_1}(x^T_1) - \alpha \right)_+$$



$\beta^*_1 \geq \beta^*_1(\epsilon)$

Distortion-free

$\epsilon$-Distorted

# Jointly Optimal Detector and Watermarking Scheme

**Optimization problem:**

$$\min_{\gamma,\, P_{X_1^T, \zeta_1^T}} \beta_1(\gamma,\, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma,\, Q_{X_1^T},\, P_{\zeta_1^T}) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T},\, Q_{X_1^T}) \leq \epsilon$$

# Jointly Optimal Detector and Watermarking Scheme

✦ **Jointly optimal detector $\gamma^*$ and watermarking scheme $P^*_{X_1^T, \zeta_1^T}$:**

**Optimization problem:**

$$\min_{\gamma, \, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma, \, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

# Jointly Optimal Detector and Watermarking Scheme

**◆ Jointly optimal detector $\gamma^*$**

**and watermarking scheme $P^*_{X_1^T, \zeta_1^T}$ :**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

for some surjective $g : \mathscr{Z}^T \to \mathcal{S} \supset \mathscr{V}^T$

**Optimization problem:**

$$\min_{\gamma,\ P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma,\ P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

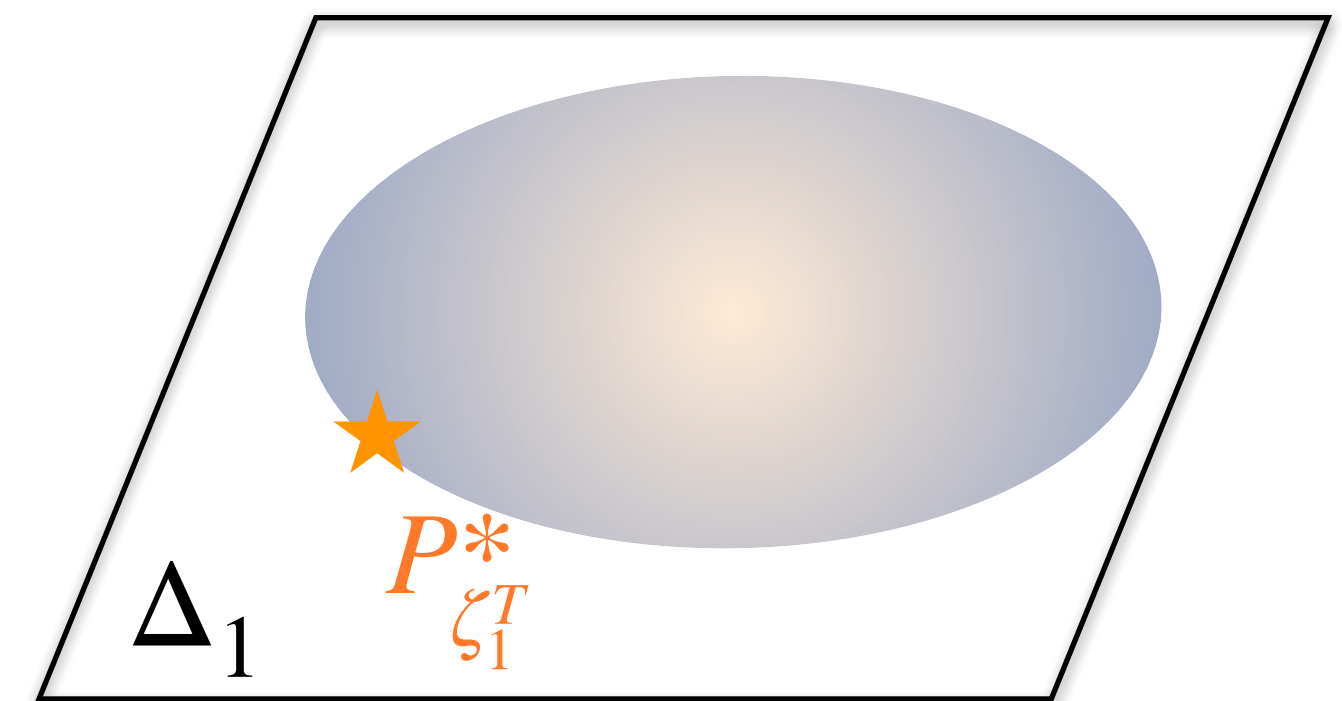# Jointly Optimal Detector and Watermarking Scheme

**Optimization problem:**

$$\min_{\gamma,\, P_{X_1^T, \zeta_1^T}} \beta_1(\gamma,\, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

✦ **Jointly optimal detector $\gamma^*$**

**and watermarking scheme $P^*_{X_1^T, \zeta_1^T}$:**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

$$\text{for some surjective } g : \mathscr{Z}^T \to \mathcal{S} \supset \mathscr{V}^T$$

$$P^*_{X_1^T, \zeta_1^T} :$$

# Jointly Optimal Detector and Watermarking Scheme

**Optimization problem:**

$$\min_{\gamma,\ P_{X_1^T,\zeta_1^T}} \quad \beta_1(\gamma,\ P_{X_1^T,\zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha \qquad (\Delta_1)$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

✦ **Jointly optimal detector $\gamma^*$ and watermarking scheme $P^*_{X_1^T,\zeta_1^T}$:**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

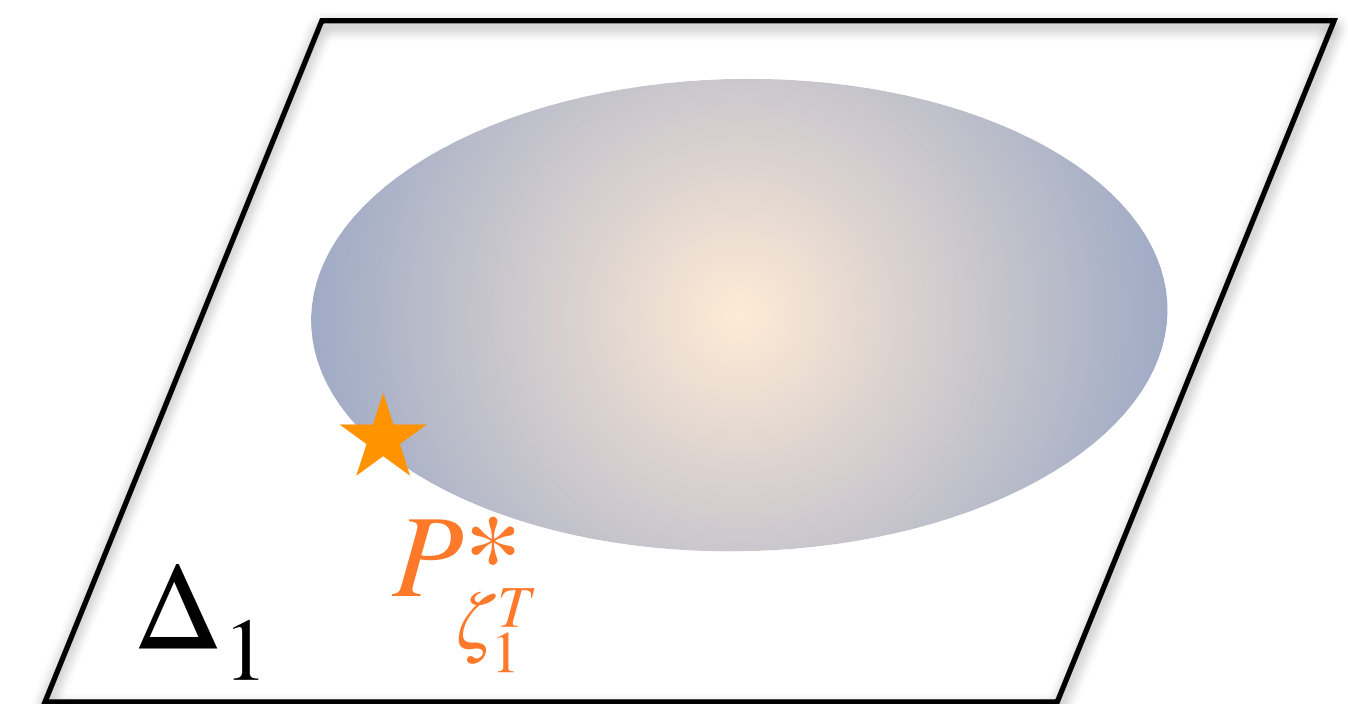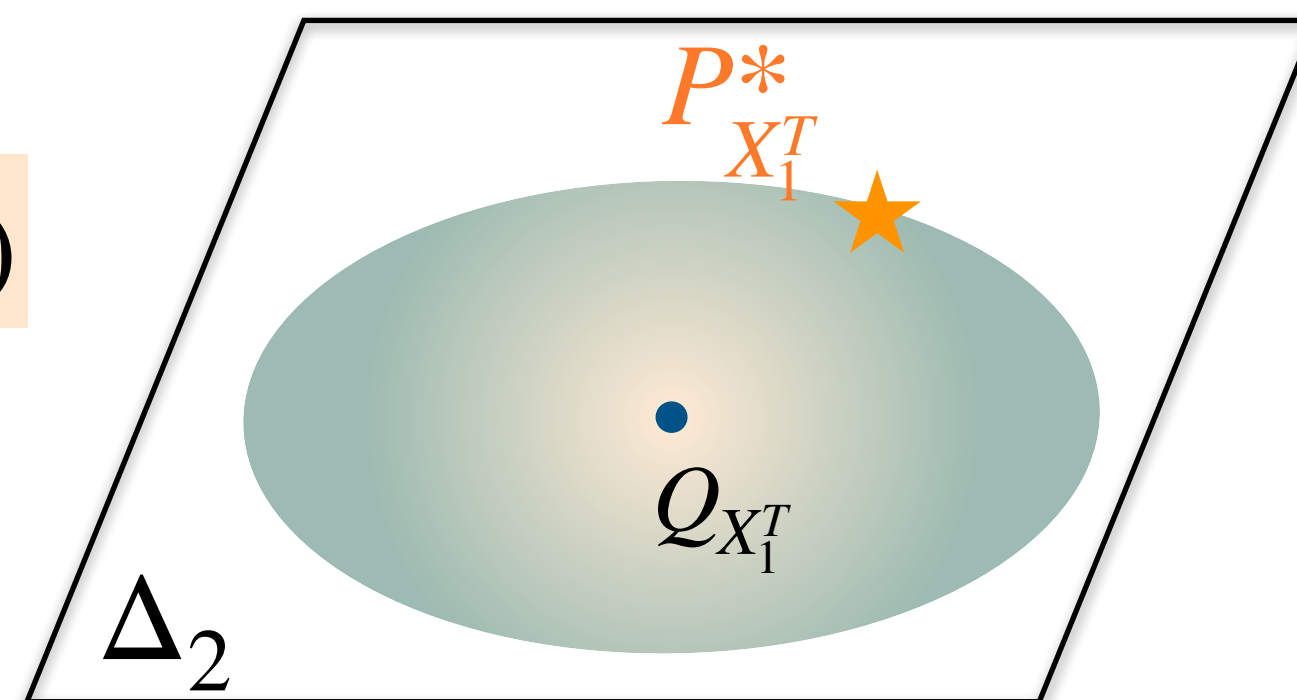$$\text{for some surjective } g : \mathcal{Z}^T \to \mathcal{S} \supset \mathcal{V}^T$$

$$P^*_{X_1^T,\zeta_1^T} :$$

# Jointly Optimal Detector and Watermarking Scheme

✦ **Jointly optimal detector $\gamma^*$ and watermarking scheme $P^*_{X_1^T, \zeta_1^T}$ :**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

for some surjective $g : \mathscr{Z}^T \to \mathcal{S} \supset \mathscr{V}^T$

**Optimization problem:**

$$\min_{\gamma, \, P_{X_1^T, \zeta_1^T}} \beta_1(\gamma, \, P_{X_1^T, \zeta_1^T})$$

s.t. $\sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$  ($\Delta_1$)

$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$  ($\Delta_2$)

$P^*_{X_1^T, \zeta_1^T}$ :



$$P^*_{X_1^T} = \arg\min_{P_{X_1^T}:\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$$

# Jointly Optimal Detector and Watermarking Scheme

✦ **Jointly optimal detector $\gamma^*$**

**and watermarking scheme $P^*_{X_1^T, \zeta_1^T}$:**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

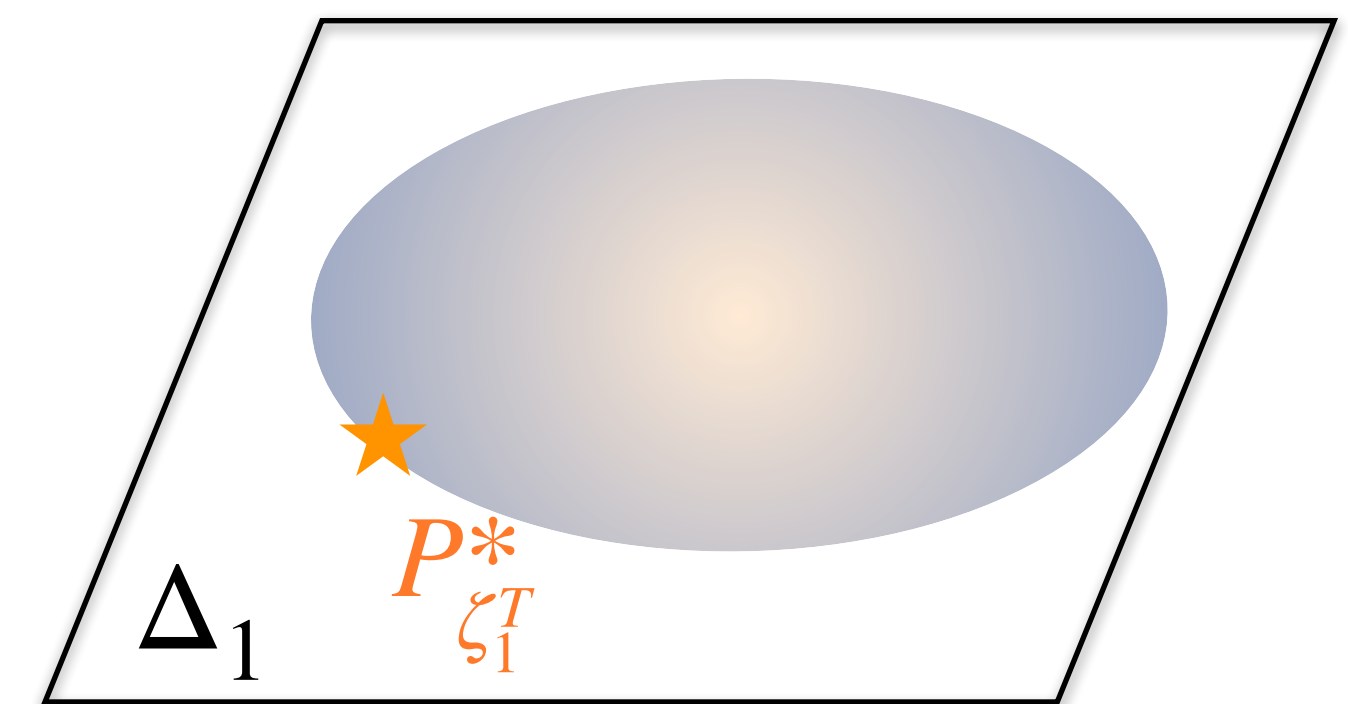for some surjective $g : \mathscr{Z}^T \to \mathcal{S} \supset \mathscr{V}^T$

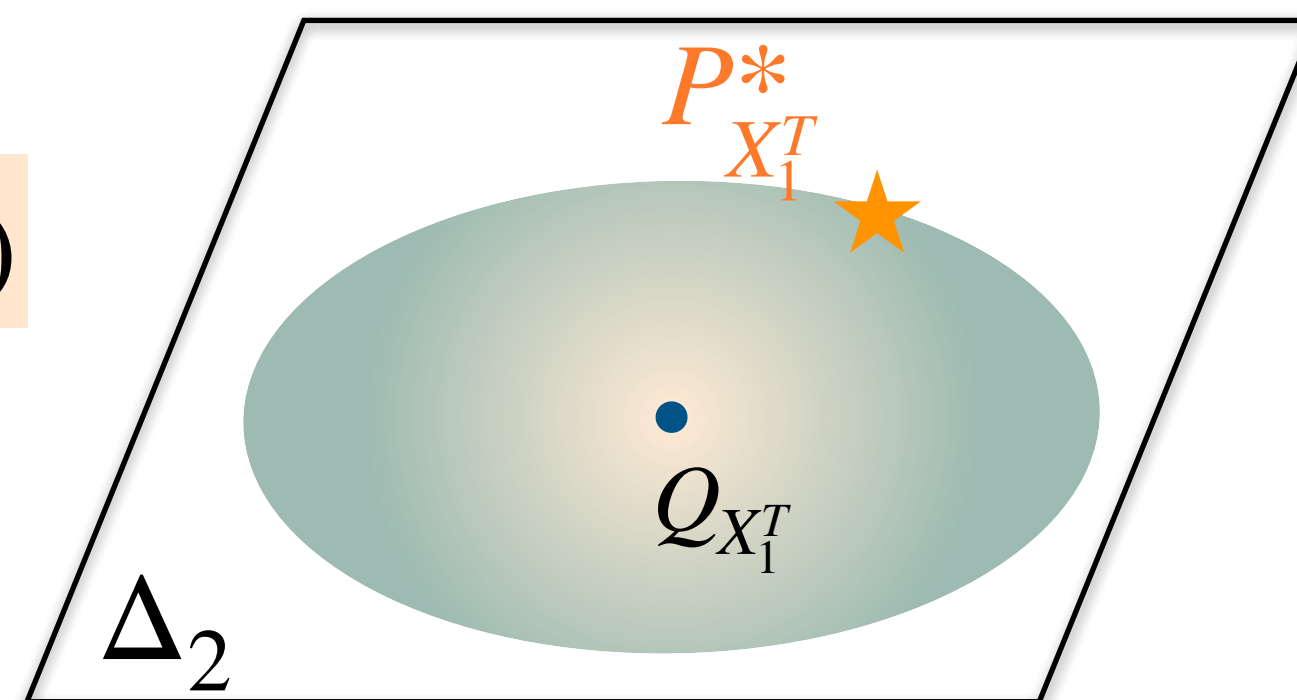**Optimization problem:**

$$\min_{\gamma,\, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma,\, P_{X_1^T, \zeta_1^T}) = \mathbb{E}_{P_{X_1^T, \zeta_1^T}}[1 - \gamma(X_1^T, \zeta_1^T)]$$

s.t. $\quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha \quad (\Delta_1)$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon \quad (\Delta_2)$$

$P^*_{X_1^T, \zeta_1^T}:$



$$P^*_{X_1^T} = \arg\min_{P_{X_1^T} : \mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$$

# Jointly Optimal Detector and Watermarking Scheme

✦ **Jointly optimal detector $\gamma^*$ and watermarking scheme $P^*_{X_1^T, \zeta_1^T}$:**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

for some surjective $g : \mathscr{Z}^T \to \mathcal{S} \supset \mathscr{V}^T$

**Optimization problem:**

$$\min_{\gamma, \, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma, \, P_{X_1^T, \zeta_1^T}) = \mathbb{E}_{P_{X_1^T, \zeta_1^T}}[1 - \gamma(X_1^T, \zeta_1^T)]$$
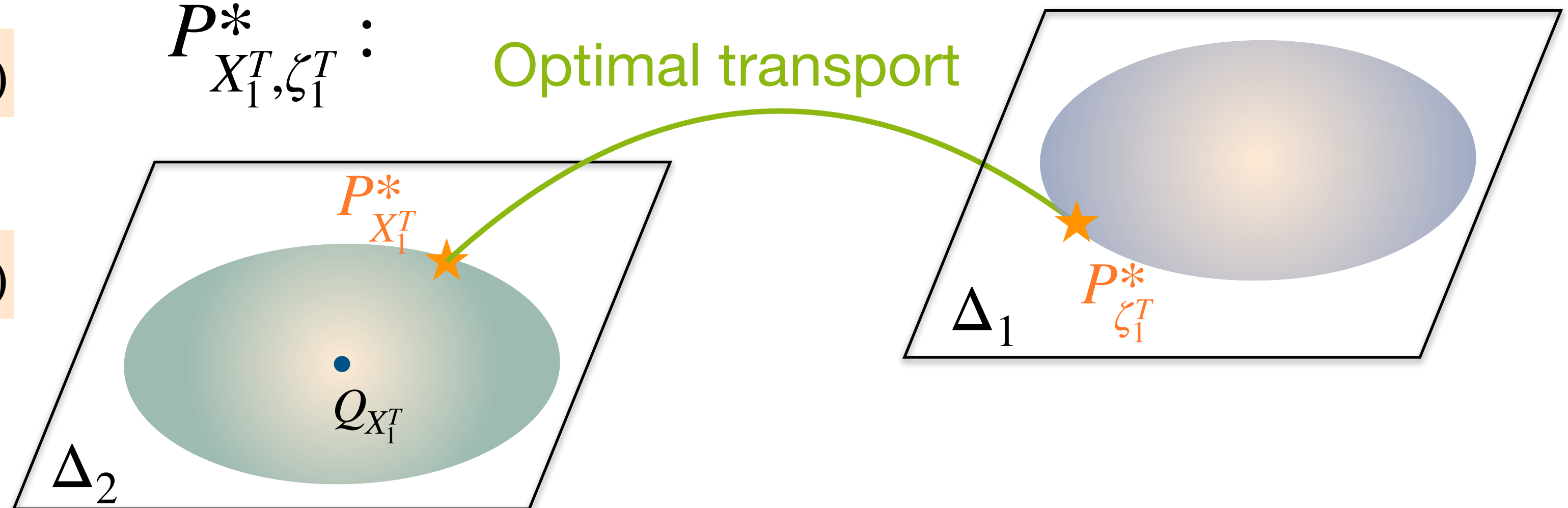
s.t. $\quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \le \alpha \quad (\Delta_1)$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \le \epsilon \quad (\Delta_2)$$

$P^*_{X_1^T, \zeta_1^T} :$

Optimal transport

$$P^*_{X_1^T} = \arg \min_{P_{X_1^T} : \mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \le \epsilon} \sum_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$$

# Jointly Optimal Detector and Watermarking Scheme

**✦ Jointly optimal detector $\gamma^*$**

**and watermarking scheme $P^*_{X_1^T, \zeta_1^T}$ :**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

for some surjective $g : \mathscr{Z}^T \to \mathcal{S} \supset \mathscr{V}^T$

**Optimization problem:**

$$\min_{\gamma,\ P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma,\ P_{X_1^T, \zeta_1^T}) = \mathbb{E}_{P_{X_1^T, \zeta_1^T}}[1 - \gamma(X_1^T, \zeta_1^T)]$$
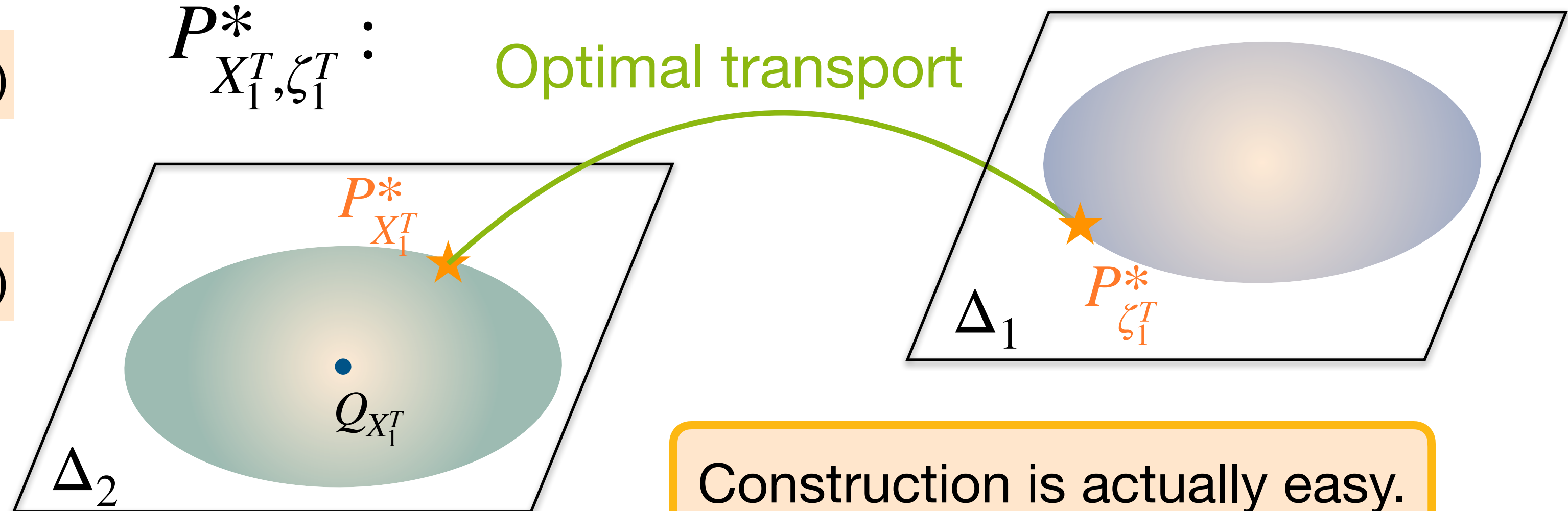
s.t. $\quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha \quad (\Delta_1)$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon \quad (\Delta_2)$$

$P^*_{X_1^T, \zeta_1^T}$ :

Optimal transport



Construction is actually easy.

$$P^*_{X_1^T} = \arg\min_{P_{X_1^T} : \mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$$

# Jointly Optimal Detector and Watermarking Scheme

**✦ Jointly optimal detector $\gamma^*$**

**and watermarking scheme $P^*_{X_1^T, \zeta_1^T}$ :**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

for some surjective $g : \mathscr{Z}^T \to \mathcal{S} \supset \mathscr{V}^T$

$$P^*_{X_1^T, \zeta_1^T} :$$

**Optimization problem:**

$$\min_{\gamma, \, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma, \, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

$$P^*_{X_1^T} = \arg \min_{P_{X_1^T} : \mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$$

# Jointly Optimal Detector and Watermarking Scheme

**Optimization problem:**

$$\min_{\gamma,\ P_{X_1^T,\zeta_1^T}} \beta_1(\gamma,\ P_{X_1^T,\zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

✦ **Jointly optimal detector $\gamma^*$**

**and watermarking scheme $P^*_{X_1^T,\zeta_1^T}$:**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

for some surjective $g : \mathscr{Z}^T \rightarrow \mathcal{S} \supset \mathscr{V}^T$

$$P^*_{X_1^T,\zeta_1^T} :$$

$$(T = 1)$$

$$P^*_{X_1^T} = \arg\min_{P_{X_1^T}:\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$$

# Jointly Optimal Detector and Watermarking Scheme

**Optimization problem:**

$$\min_{\gamma,\, P_{X_1^T, \zeta_1^T}} \beta_1(\gamma,\, P_{X_1^T, \zeta_1^T})$$

s.t. $\displaystyle\sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \le \alpha$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \le \epsilon$$

$$P_{X_1^T}^* = \arg\min_{P_{X_1^T}:\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \le \epsilon} \sum_{x_1^T} (P_{X_1^T}(x_1^T) - \alpha)_+$$

✦ **Jointly optimal detector** $\gamma^*$ **and watermarking scheme** $P_{X_1^T, \zeta_1^T}^*$ **:**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

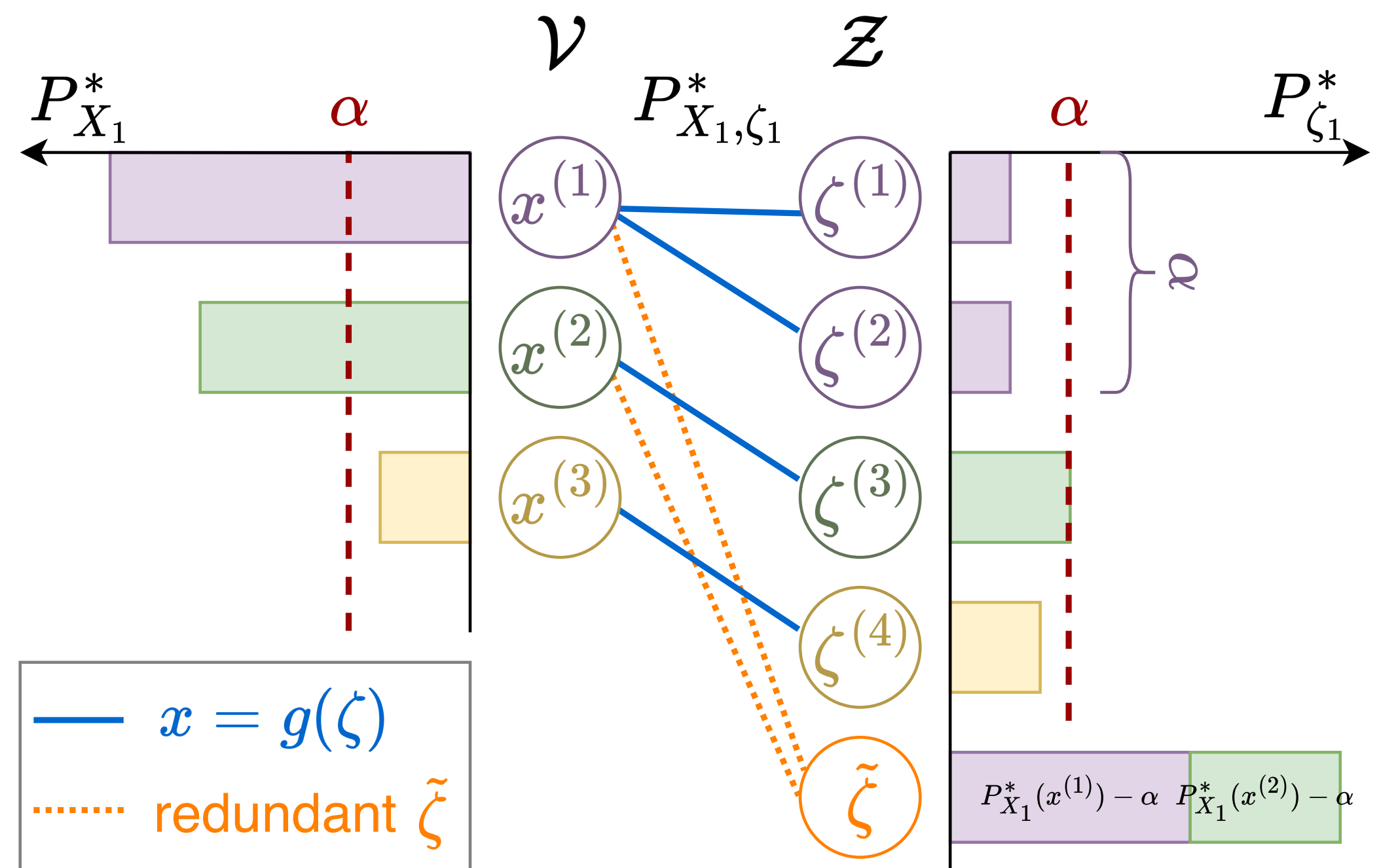for some surjective $g : \mathscr{Z}^T \to \mathcal{S} \supset \mathscr{V}^T$

$P_{X_1^T, \zeta_1^T}^* :$

$(T = 1)$



$x = g(\zeta)$

redundant $\tilde{\zeta}$

# Jointly Optimal Detector and Watermarking Scheme

✦ **Jointly optimal detector $\gamma^*$**

**and watermarking scheme $P^*_{X_1^T, \zeta_1^T}$:**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

for some surjective $g : \mathscr{Z}^T \to \mathcal{S} \supset \mathscr{V}^T$

$$P^*_{X_1^T, \zeta_1^T} :$$

**Optimization problem:**

$$\min_{\gamma,\, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

# Jointly Optimal Detector and Watermarking Scheme

✦ **Jointly optimal detector $\gamma*$**

**and watermarking scheme $P^*_{X_1^T, \zeta_1^T}$ :**

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

for some surjective $g : \mathscr{Z}^T \to \mathcal{S} \supset \mathscr{V}^T$

$P^*_{X_1^T, \zeta_1^T} :$

$P^*_{\zeta_1^T}$ **adaptive** to original LLM

predicted distribution $Q_{X_1^T}$

**Optimization problem:**

$$\min_{\gamma, \, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

# Jointly Optimal Detector and Watermarking Scheme

**✦ Jointly optimal detector $\gamma^*$**

**and watermarking scheme $P^*_{X_1^T, \zeta_1^T}$ :**

**Optimization problem:**

$$\min_{\gamma,\ P_{X_1^T, \zeta_1^T}} \beta_1(\gamma, P_{X_1^T, \zeta_1^T})$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}) \leq \alpha$$

$$D(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

$$\gamma^* = \mathbf{1}\{X_1^T = g(\zeta_1^T)\}$$

for some surjective $g : \mathscr{Z}^T \to \mathcal{S} \supset \mathscr{V}^T$

$P^*_{X_1^T, \zeta_1^T} :$ 

$P^*_{\zeta_1^T}$ **adaptive** to original LLM predicted distribution $Q_{X_1^T}$

Unlike existing watermarking methods

# Sequence-Level Optimal to Token-Level Optimal

# Sequence-Level Optimal to Token-Level Optimal

- Previous optimal result holds for fixed $T \Rightarrow$ unable to implement dynamically

# Sequence-Level Optimal to Token-Level Optimal

- Previous optimal result holds for fixed $T \Rightarrow$ unable to implement dynamically

- Solution: implement it token by token

# Sequence-Level Optimal to Token-Level Optimal

- Previous optimal result holds for fixed $T \Rightarrow$ unable to implement dynamically

- Solution: implement it token by token

  **Detector:**

# Sequence-Level Optimal to Token-Level Optimal

- Previous optimal result holds for fixed $T \Rightarrow$ unable to implement dynamically

- Solution: implement it token by token

**Detector:**

$$\gamma_{tk} = \mathbf{1}\left\{\frac{1}{T}\sum_{t=1}^{T}\mathbf{1}\{X_t = g(\zeta_t)\} \geq \lambda\right\} \quad \text{for some surjective } g : \mathscr{Z} \to \mathcal{S} \supset \mathscr{V}$$

# Sequence-Level Optimal to Token-Level Optimal

- Previous optimal result holds for fixed $T \Rightarrow$ unable to implement dynamically

- Solution: implement it token by token

**Detector:**

$$\gamma_{tk} = \mathbf{1}\left\{\frac{1}{T}\sum_{t=1}^{T}\mathbf{1}\{X_t = g(\zeta_t)\} \geq \lambda\right\} \quad \text{for some surjective } g : \mathcal{Z} \to \mathcal{S} \supset \mathcal{V}$$

**Watermarking scheme:**

$$P^*_{X_1^T, \zeta_1^T} \to P^*_{X_t, \zeta_t | X_1^{t-1}, \zeta_1^{t-1}}, \forall t = 1, 2, \ldots, T$$

# Sequence-Level Optimal to Token-Level Optimal

- Previous optimal result holds for fixed $T \Rightarrow$ unable to implement dynamically

- Solution: implement it token by token

**Detector:**

$$\gamma_{tk} = \mathbf{1}\left\{ \frac{1}{T} \sum_{t=1}^{T} \mathbf{1}\{X_t = g(\zeta_t)\} \geq \lambda \right\} \quad \text{for some surjective } g : \mathscr{X} \to \mathscr{S} \supset \mathscr{V}$$
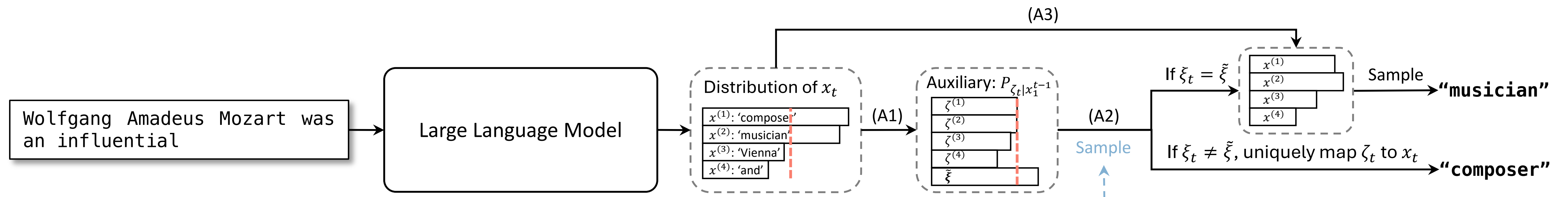
**Watermarking scheme:**

$$P^*_{X_1^T, \zeta_1^T} \to P^*_{X_t, \zeta_t | X_1^{t-1}, \zeta_1^{t-1}}, \forall t = 1, 2, \ldots, T$$

dependent

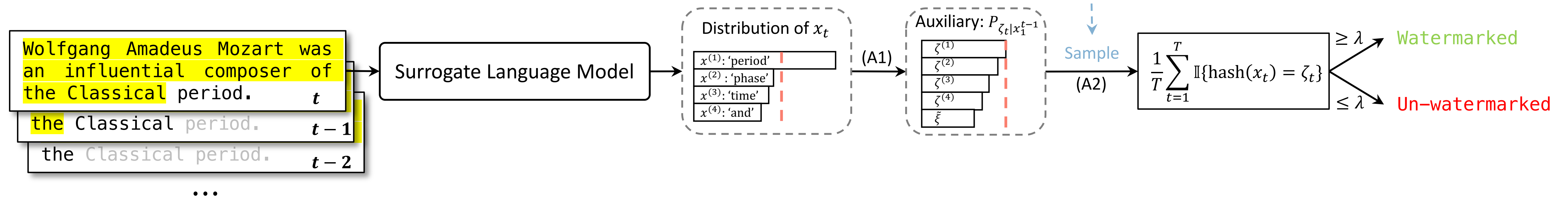Token-level false alarm rate $\eta \xrightarrow{controls}$ Sequence-level false alarm $\alpha$

# DAWA: Distribution-Adaptive Watermarking Algorithm ($\epsilon = 0$, distortion-free)
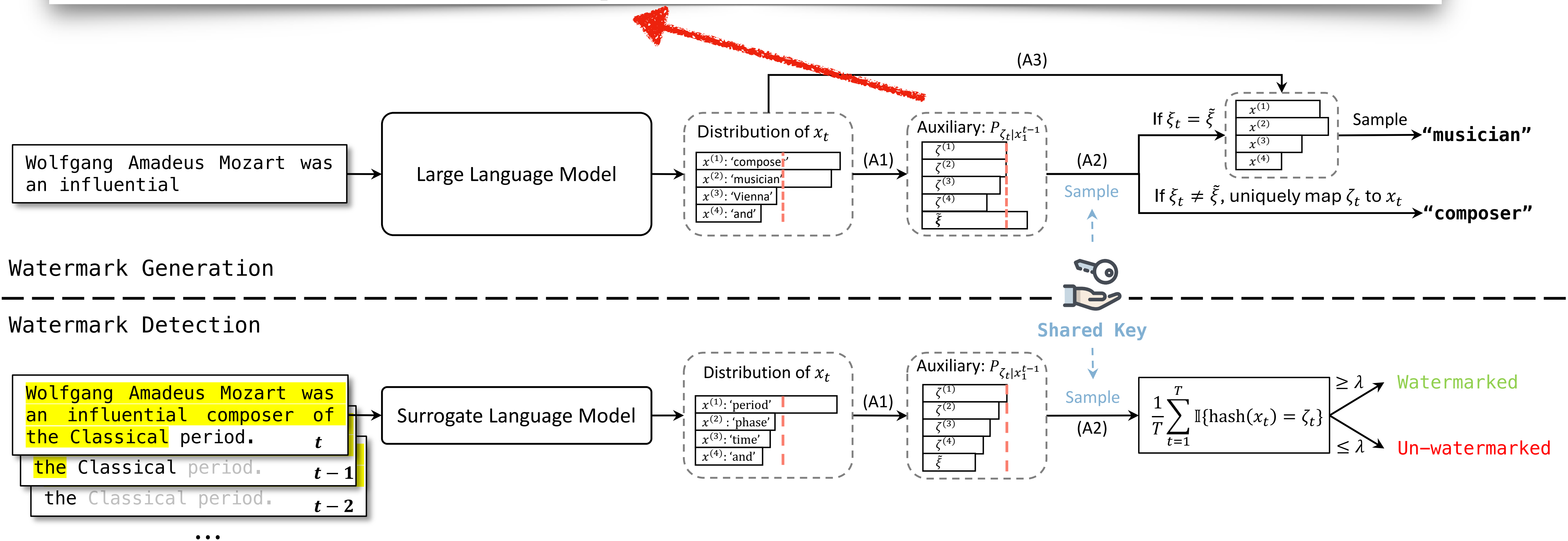
# DAWA: Distribution-Adaptive Watermarking Algorithm ($\epsilon = 0$, distortion-free)

At each time $t$, construct $P^*_{\zeta_t|X_1^t}$ from the LLM predicted distribution $Q_{X_t|X_1^{t-1}}$



Watermark Generation

Watermark Detection

Shared Key

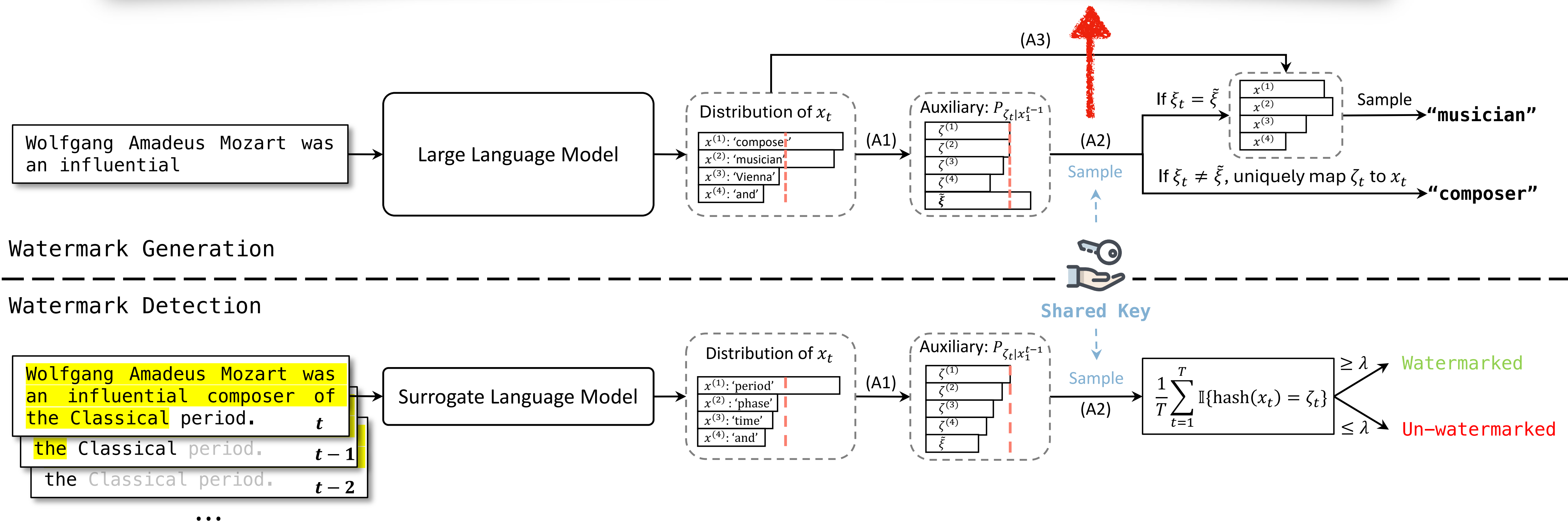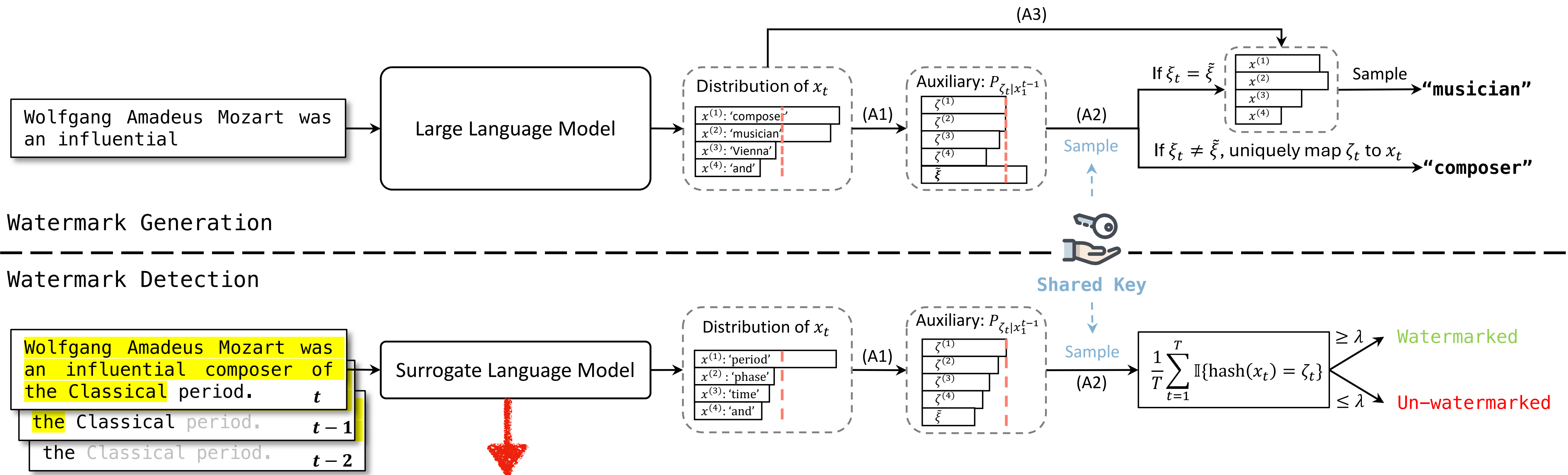# DAWA: Distribution-Adaptive Watermarking Algorithm ($\epsilon = 0$, distortion-free)

Sample $\zeta_t$ using Gumbel max trick: $\zeta_t \leftarrow \arg\max_{\zeta} \log P^*_{\zeta_t|x_1^t}(\zeta) + G_{\zeta,t}$
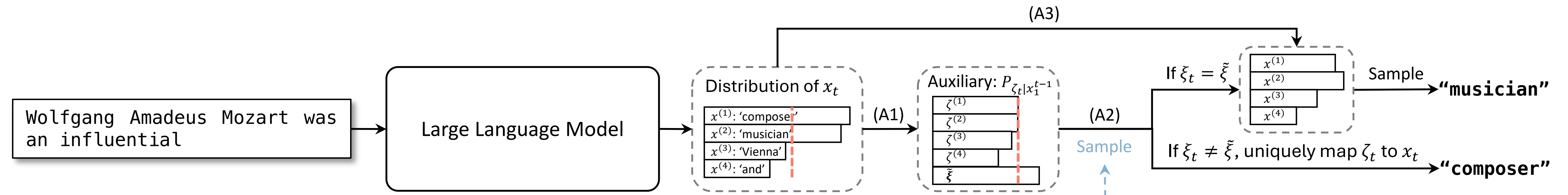
(A3)

Wolfgang Amadeus Mozart was an influential

Large Language Model

Distribution of $x_t$

$x^{(1)}$: 'composer'
$x^{(2)}$: 'musician'
$x^{(3)}$: 'Vienna'
$x^{(4)}$: 'and'

(A1)

Auxiliary: $P_{\zeta_t|x_1^{t-1}}$

$\zeta^{(1)}$
$\zeta^{(2)}$
$\zeta^{(3)}$
$\zeta^{(4)}$
$\tilde{\xi}$

(A2)

Sample

If $\xi_t = \tilde{\xi}$

$x^{(1)}$
$x^{(2)}$
$x^{(3)}$
$x^{(4)}$

Sample → **"musician"**

If $\xi_t \neq \tilde{\xi}$, uniquely map $\zeta_t$ to $x_t$

→ **"composer"**

Watermark Generation

**Shared Key**

Watermark Detection

Wolfgang Amadeus Mozart was an influential composer of the Classical period. $t$

the Classical period. $t-1$

the Classical period. $t-2$

...

Surrogate Language Model

Distribution of $x_t$

$x^{(1)}$: 'period'
$x^{(2)}$: 'phase'
$x^{(3)}$: 'time'
$x^{(4)}$: 'and'

(A1)

Auxiliary: $P_{\zeta_t|x_1^{t-1}}$

$\zeta^{(1)}$
$\zeta^{(2)}$
$\zeta^{(3)}$
$\zeta^{(4)}$
$\tilde{\xi}$

Sample

(A2)

$\frac{1}{T}\sum_{t=1}^{T} \mathbb{I}\{\text{hash}(x_t) = \zeta_t\}$

$\geq \lambda$ → Watermarked

$\leq \lambda$ → Un-watermarked

# DAWA: Distribution-Adaptive Watermarking Algorithm ($\epsilon = 0$, distortion-free)



Watermark Generation

Watermark Detection

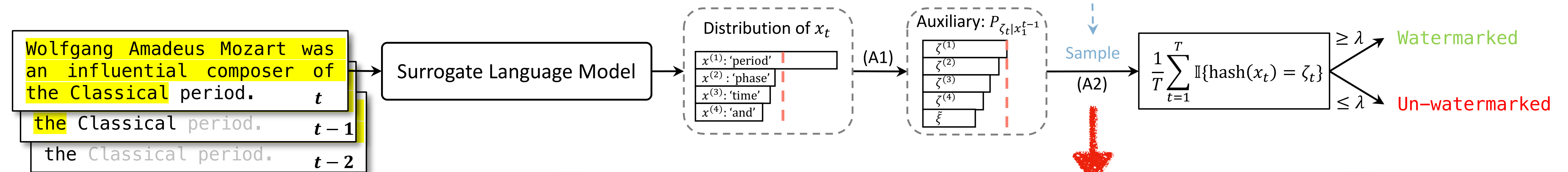Approximate distribution of $X_t$ so as to construct $\tilde{P}_{\zeta_t|x_1^t}$

# DAWA: Distribution-Adaptive Watermarking Algorithm ($\epsilon = 0$, distortion-free)



**Watermark Generation**

**Watermark Detection**

Sample $\zeta_t$ using Gumbel max trick: $\zeta_t \leftarrow \arg\max_{\zeta} \log \tilde{P}_{\zeta_t | x_1^t}(\zeta) + G_{\zeta,t}$
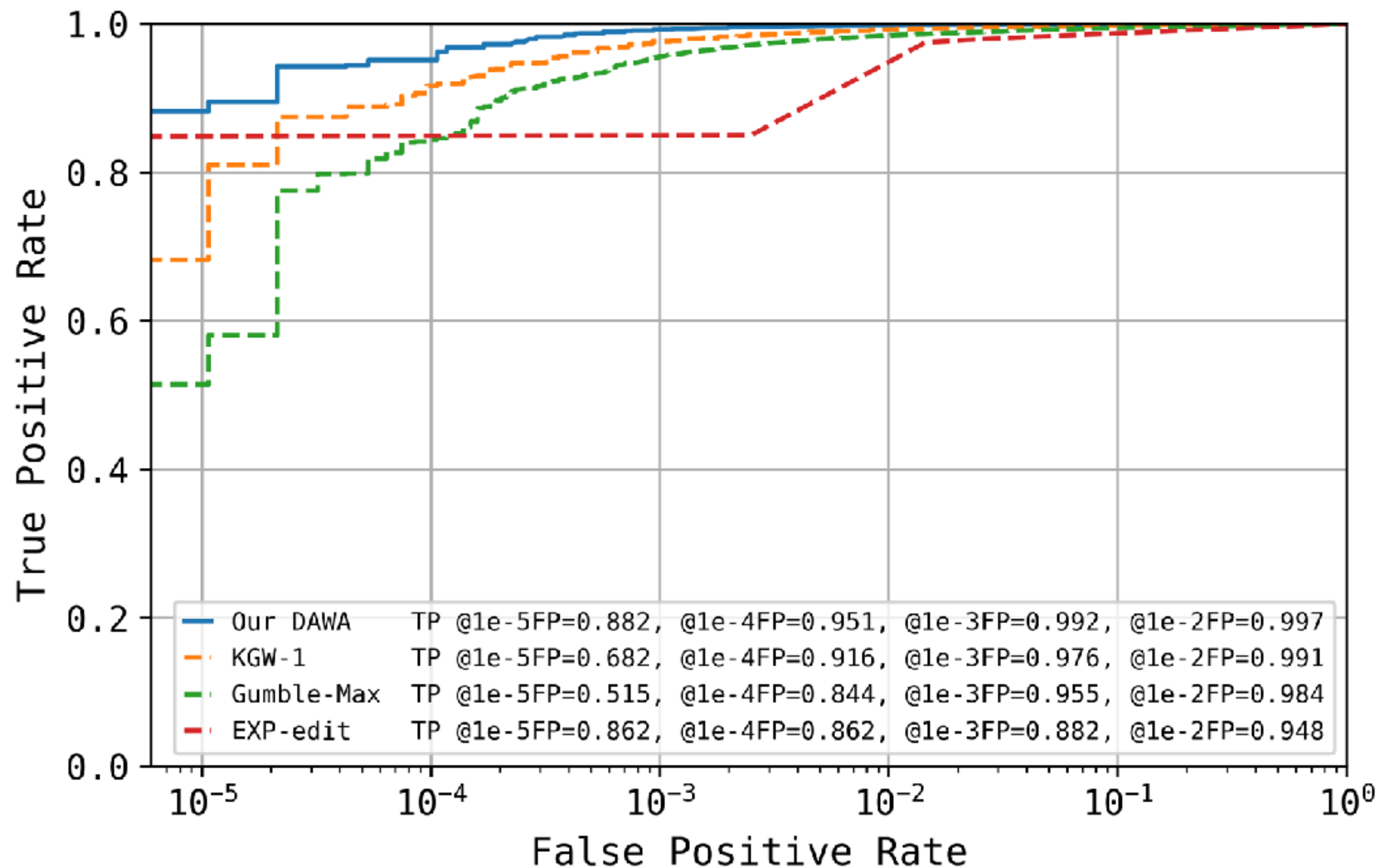
# Experimental Result

**DAWA** (**D**istribution-**A**daptive **W**atermarking **A**lgorithm)

# Experimental Result

**DAWA** (**D**istribution-**A**daptive **W**atermarking **A**lgorithm)

Fast and Accurate

# Experimental Result

**DAWA** (**D**istribution-**A**daptive **W**atermarking **A**lgorithm)

**Fast and Accurate**

**Text quality high**

| Methods | Human | KGW-1 | EXP-Edit | Gumbel-Max | Ours |
|---|---|---|---|---|---|
| BLEU Score | 0.219 | 0.158 | 0.203 | 0.210 | 0.214 |
| Avg Perplexity | 8.846 | 14.327 | 12.186 | 11.732 | 6.495 |

# Future Work: With Text Modifications?

# Future Work: With Text Modifications?

- $f : \mathcal{V}^T \to [K]$: a function that maps a sequence of tokens $X_1^T$ to a finite latent space $[K]$, e.g., a semantic mapping

# Future Work: With Text Modifications?

- $f : \mathcal{V}^T \to [K]$: a function that maps a sequence of tokens $X_1^T$ to a finite latent space $[K]$, e.g., a semantic mapping

- $\mathcal{B}_f(x_1^T) = \{\tilde{x}_1^T \in \mathcal{V}^T : f(\tilde{x}_1^T) = f(x_1^T)\}$ be an equivalence class containing $x_1^T$

# Future Work: With Text Modifications?

- $f : \mathscr{V}^T \to [K]$: a function that maps a sequence of tokens $X_1^T$ to a finite latent space $[K]$, e.g., a semantic mapping

- $\mathscr{B}_f(x_1^T) = \{\tilde{x}_1^T \in \mathscr{V}^T : f(\tilde{x}_1^T) = f(x_1^T)\}$ be an equivalence class containing $x_1^T$

- Text modification: $x_1^T$ can be modified as any text within $\mathscr{B}_f(x_1^T)$

# Future Work: With Text Modifications?

- $f : \mathcal{V}^T \to [K]$: a function that maps a sequence of tokens $X_1^T$ to a finite latent space $[K]$, e.g., a semantic mapping

- $\mathscr{B}_f(x_1^T) = \{\tilde{x}_1^T \in \mathcal{V}^T : f(\tilde{x}_1^T) = f(x_1^T)\}$ be an equivalence class containing $x_1^T$

- Text modification: $x_1^T$ can be modified as any text within $\mathscr{B}_f(x_1^T)$

- $f$-robust Type-I and Type-II errors:

# Future Work: With Text Modifications?

- $f : \mathcal{V}^T \to [K]$: a function that maps a sequence of tokens $X_1^T$ to a finite latent space $[K]$, e.g., a semantic mapping

- $\mathcal{B}_f(x_1^T) = \{\tilde{x}_1^T \in \mathcal{V}^T : f(\tilde{x}_1^T) = f(x_1^T)\}$ be an equivalence class containing $x_1^T$

- Text modification: $x_1^T$ can be modified as any text within $\mathcal{B}_f(x_1^T)$

- $f$-robust Type-I and Type-II errors:

$$\beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}, f) := \mathbb{E}_{Q_{X_1^T} \otimes P_{\zeta_1^T}} \left[ \sup_{\tilde{x}_1^T \in \mathcal{B}_f(X_1^T)} \mathbf{1}\{\gamma(\tilde{x}_1^T, \zeta_1^T) = 1\} \right]$$

# Future Work: With Text Modifications?

- $f : \mathscr{V}^T \to [K]$: a function that maps a sequence of tokens $X_1^T$ to a finite latent space $[K]$, e.g., a semantic mapping

- $\mathscr{B}_f(x_1^T) = \{\tilde{x}_1^T \in \mathscr{V}^T : f(\tilde{x}_1^T) = f(x_1^T)\}$ be an equivalence class containing $x_1^T$

- Text modification: $x_1^T$ can be modified as any text within $\mathscr{B}_f(x_1^T)$

- $f$-robust Type-I and Type-II errors:

$$\beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}, f) := \mathbb{E}_{Q_{X_1^T} \otimes P_{\zeta_1^T}} \left[ \sup_{\tilde{x}_1^T \in \mathscr{B}_f(X_1^T)} \mathbf{1}\{\gamma(\tilde{x}_1^T, \zeta_1^T) = 1\} \right]$$

$$\beta_1(\gamma, P_{X_1^T, \zeta_1^T}, f) := \mathbb{E}_{P_{X_1^T, \zeta_1^T}} \left[ \sup_{\tilde{x}_1^T \in \mathscr{B}_f(X_1^T)} \mathbf{1}\{\gamma(\tilde{x}_1^T, \zeta_1^T) = 0\} \right]$$

# **Robustness Against Text Modifications**

**<u>Optimization problem:</u>**

$$\min_{\gamma,\, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma,\, P_{X_1^T, \zeta_1^T}, f)$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}, f) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

# Robustness Against Text Modifications

**<u>Optimization problem:</u>**

$$\min_{\gamma,\ P_{X_1^T, \zeta_1^T}} \beta_1(\gamma,\ P_{X_1^T, \zeta_1^T}, f)$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}, f) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

✦ **Minimum $f$-robust Type-II error:**

# Robustness Against Text Modifications

**<u>Optimization problem:</u>**

$$\min_{\gamma,\ P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma,\ P_{X_1^T, \zeta_1^T}, f)$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}, f) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

✦ **Minimum $f$-robust Type-II error:**

$$\beta_1^*(Q_{X_1^T}, \alpha, \epsilon, f)$$

$$= \min_{P_{X_1^T} : \mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum_{k \in [K]} \left( \left( \sum_{x_1^T : f(x_1^T) = k} P_{X_1^T}(x_1^T) \right) - \alpha \right)_+$$

# Robustness Against Text Modifications

**Optimization problem:**

$$\min_{\gamma, \, P_{X_1^T, \zeta_1^T}} \quad \beta_1(\gamma, \, P_{X_1^T, \zeta_1^T}, f)$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma, Q_{X_1^T}, P_{\zeta_1^T}, f) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon$$

✦ **Minimum $f$-robust Type-II error:**

$$\beta_1^*(Q_{X_1^T}, \alpha, \epsilon, f)$$

$$= \min_{P_{X_1^T} : \mathsf{D}(P_{X_1^T}, Q_{X_1^T}) \leq \epsilon} \sum_{k \in [K]} \left( \left( \sum_{x_1^T : f(x_1^T) = k} P_{X_1^T}(x_1^T) \right) - \alpha \right)_+$$

Higher than the minimum Type-II error without considering robustness

# Robustness Against Text Modifications

**<u>Optimization problem:</u>**

$$\min_{\gamma,\ P_{X_1^T,\zeta_1^T}} \quad \beta_1(\gamma,\ P_{X_1^T,\zeta_1^T},f)$$

$$\text{s.t.} \quad \sup_{Q_{X_1^T}} \beta_0(\gamma,Q_{X_1^T},P_{\zeta_1^T},f) \leq \alpha$$

$$\mathsf{D}(P_{X_1^T},Q_{X_1^T}) \leq \epsilon$$

✦ **Minimum $f$-robust Type-II error:**

$$\beta_1^*(Q_{X_1^T},\alpha,\epsilon,f)$$

$$= \min_{P_{X_1^T}:\mathsf{D}(P_{X_1^T},Q_{X_1^T})\leq\epsilon} \sum_{k\in[K]} \left(\left(\sum_{x_1^T:f(x_1^T)=k} P_{X_1^T}(x_1^T)\right) - \alpha\right)_+$$

Higher than the minimum Type-II error without considering robustness

✦ **Optimal watermarking scheme:**

add signal $\zeta_1^T$ to $P_{f(X_1^T)}$, e.g., in the semantic space