



Obliviator Reveals the Cost of *Nonlinear Guardedness* in Concept Erasure



Ramin Akbari*



Milad Afshari*



Vishnu Naresh Boddeti

*Equal Contribution

Why Do We Need Erasure ?

She took a Ph.D. in Earth and Atmospheric Sciences from the Georgia Institute of Technology. Her current research involves a study of the central Pacific climate and ENSO variability over the past 6000 years.



X



Information of Interest : Profession

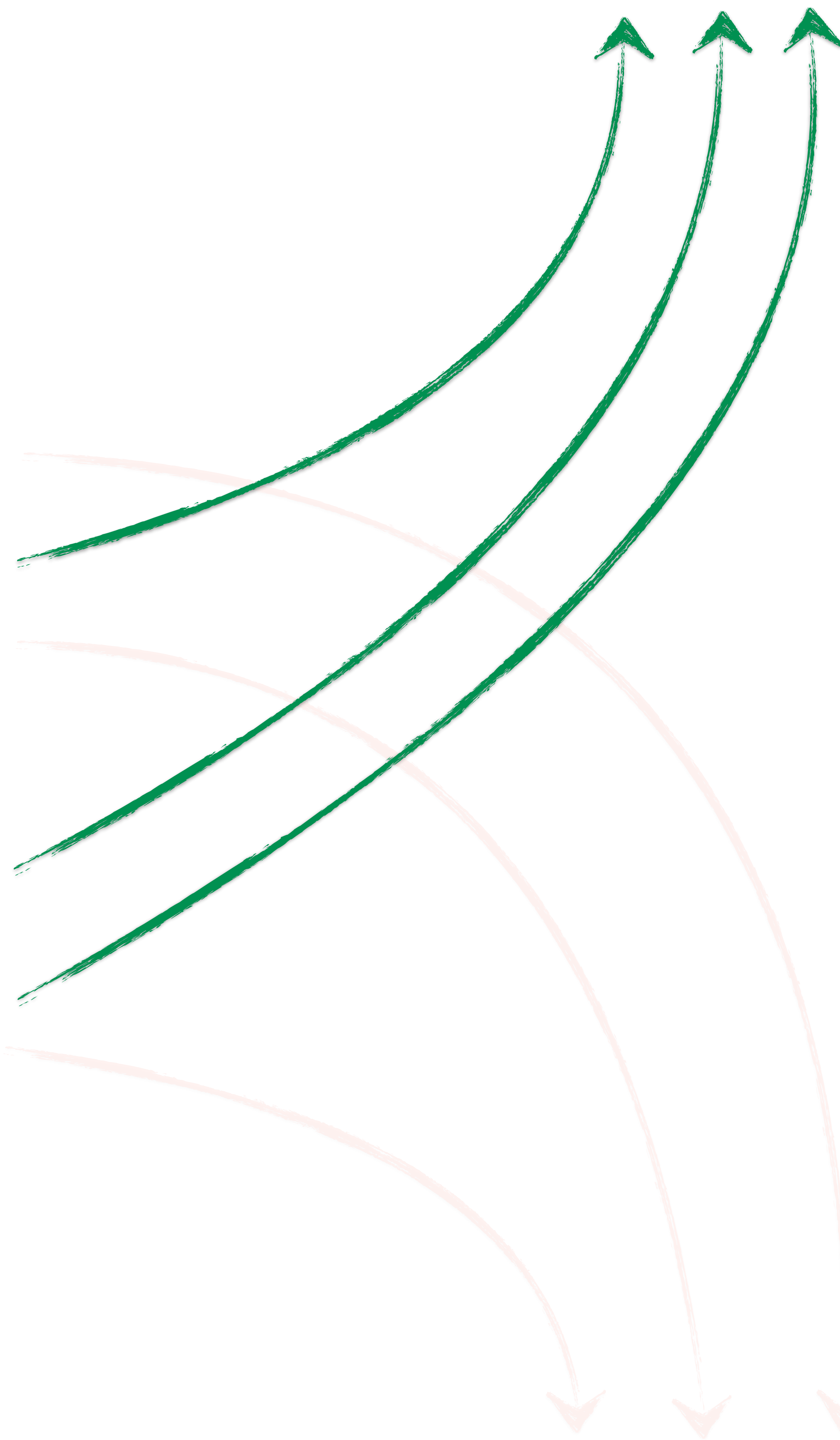
Why Do We Need Erasure ?

She took a Ph.D. in Earth and Atmospheric Sciences from the Georgia Institute of Technology. Her current research involves a study of the central Pacific climate and ENSO variability over the past 6000 years.



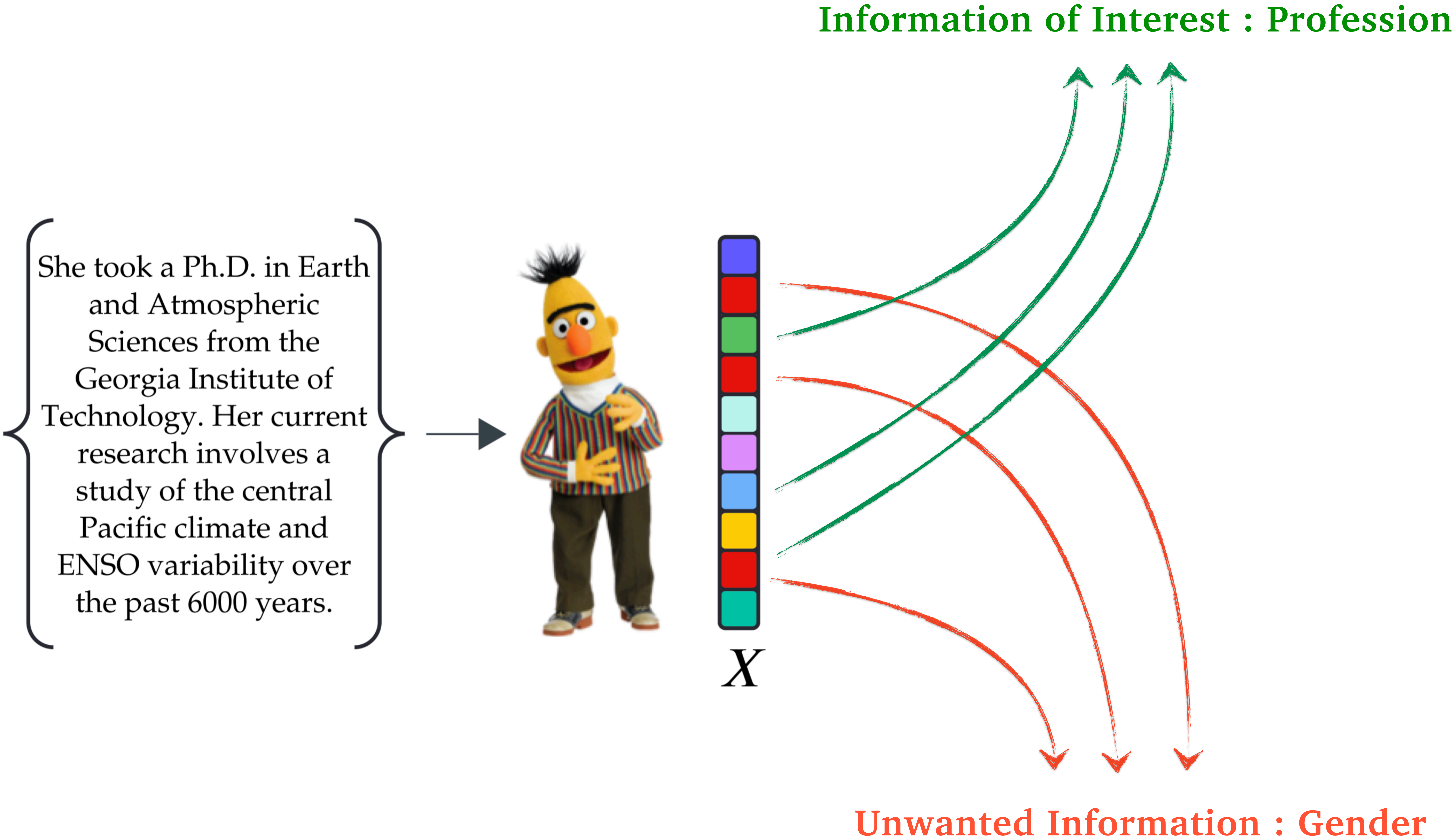
X

Information of Interest : Profession

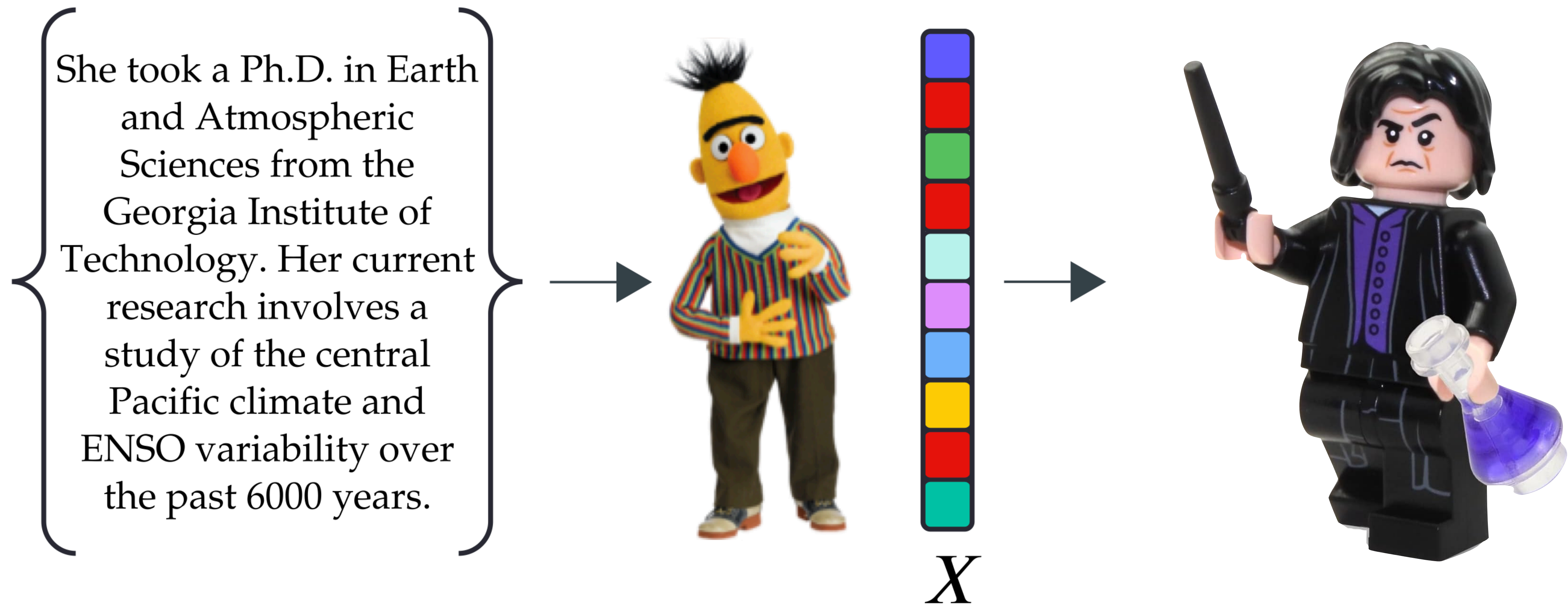


Unwanted Information : Gender

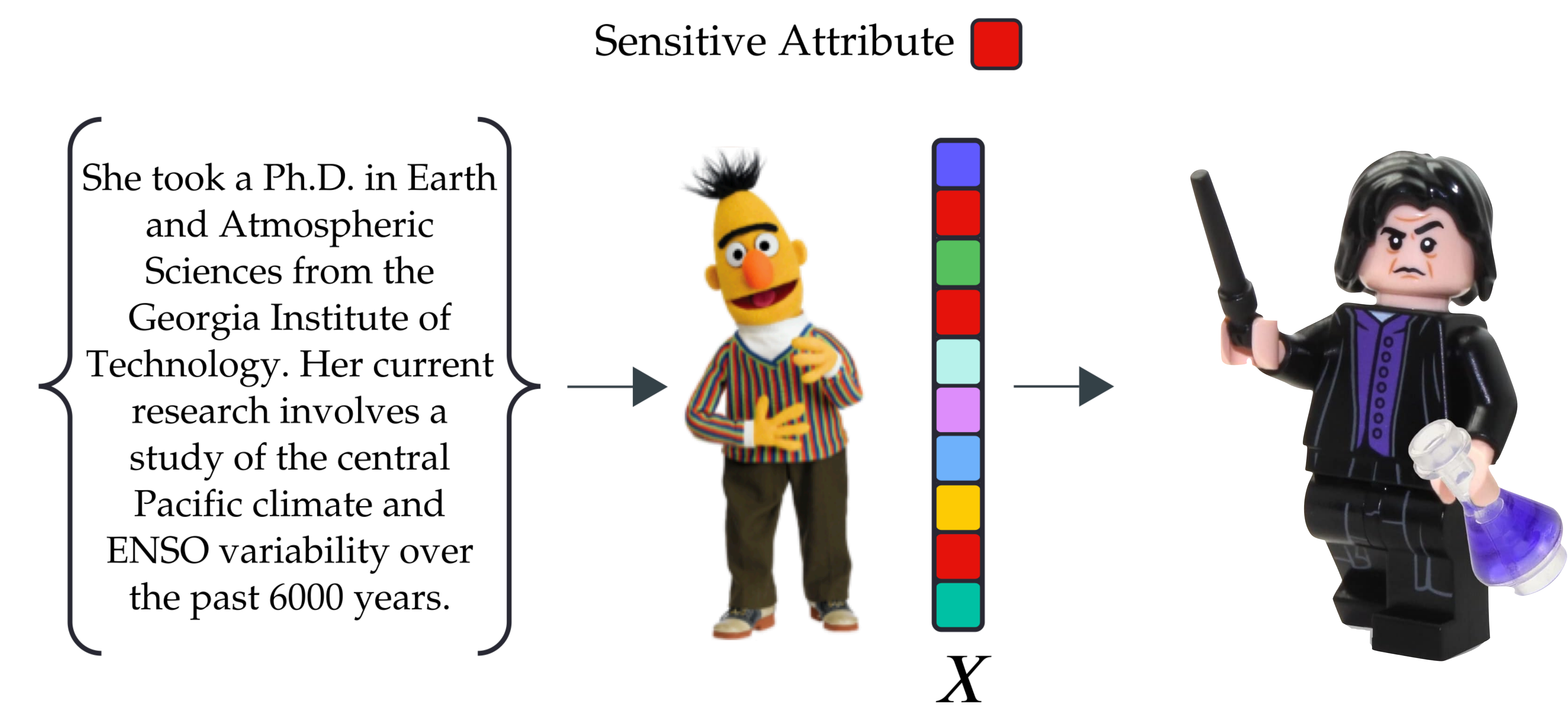
Why Do We Need Erasure ?



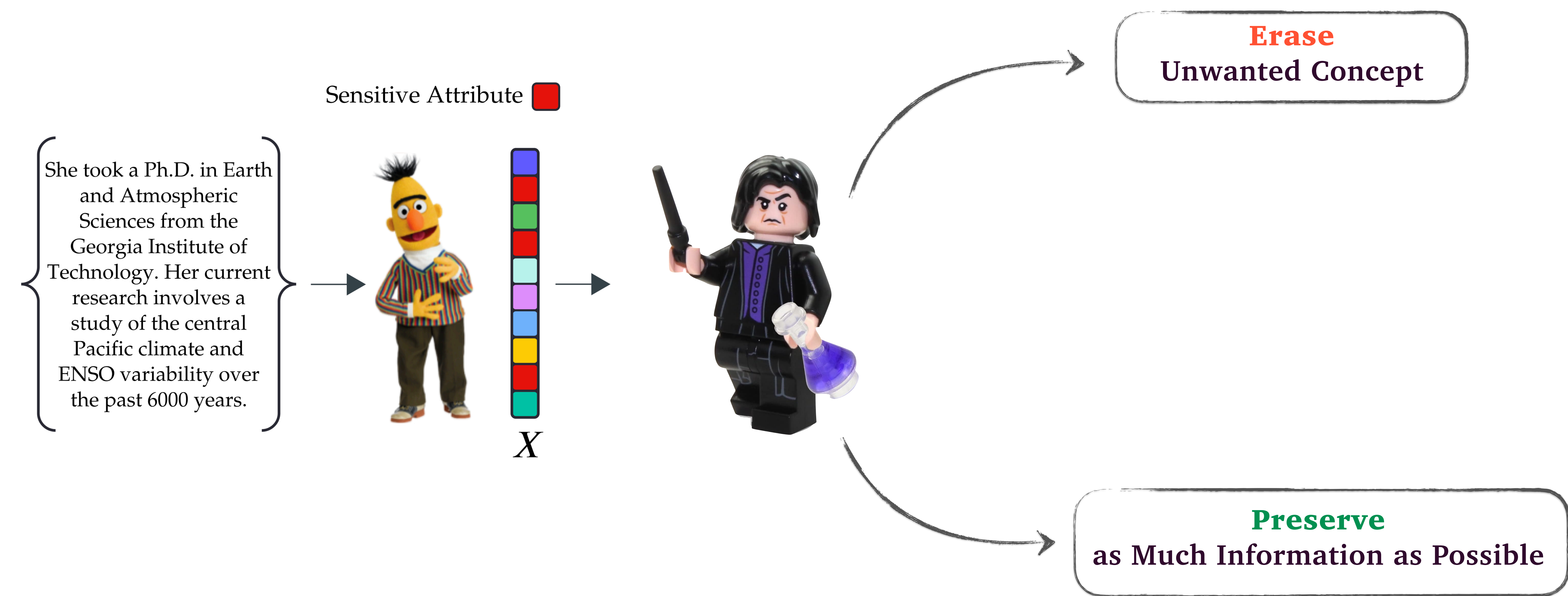
What Is The **Ultimate** Goal of Concept Erasure ?



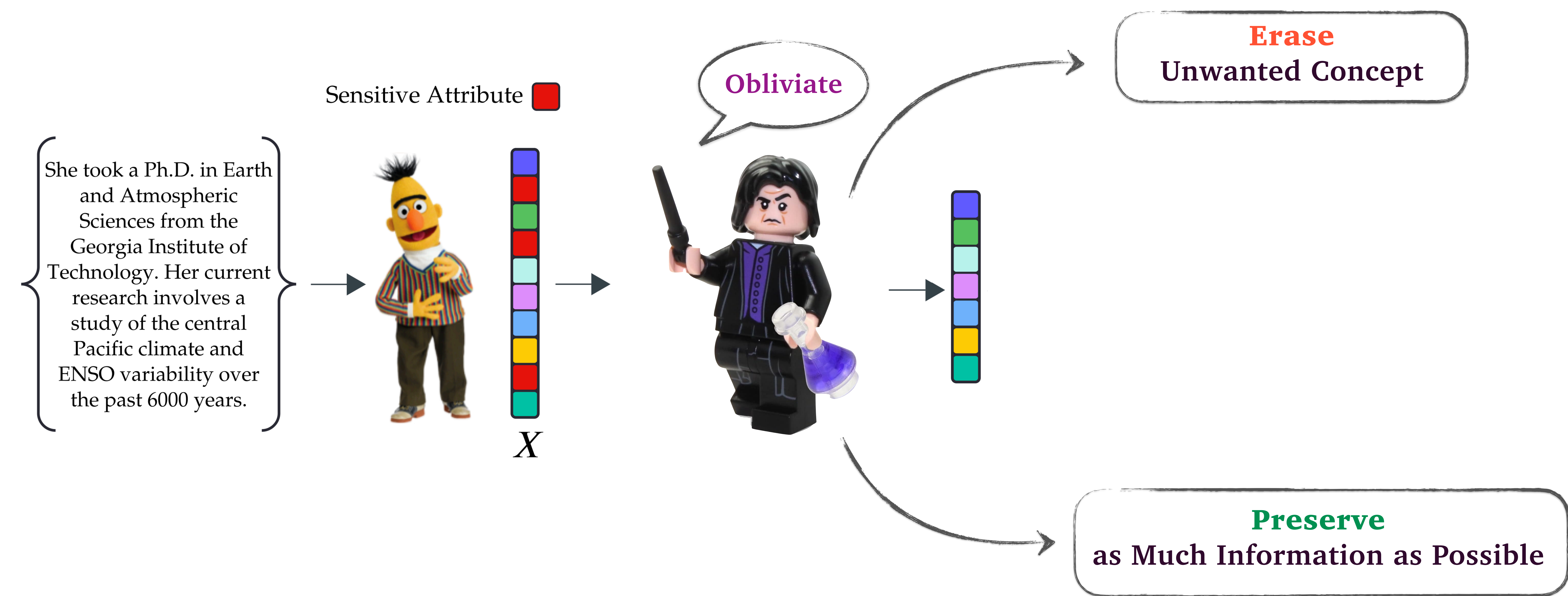
What Is The **Ultimate** Goal of Concept Erasure ?



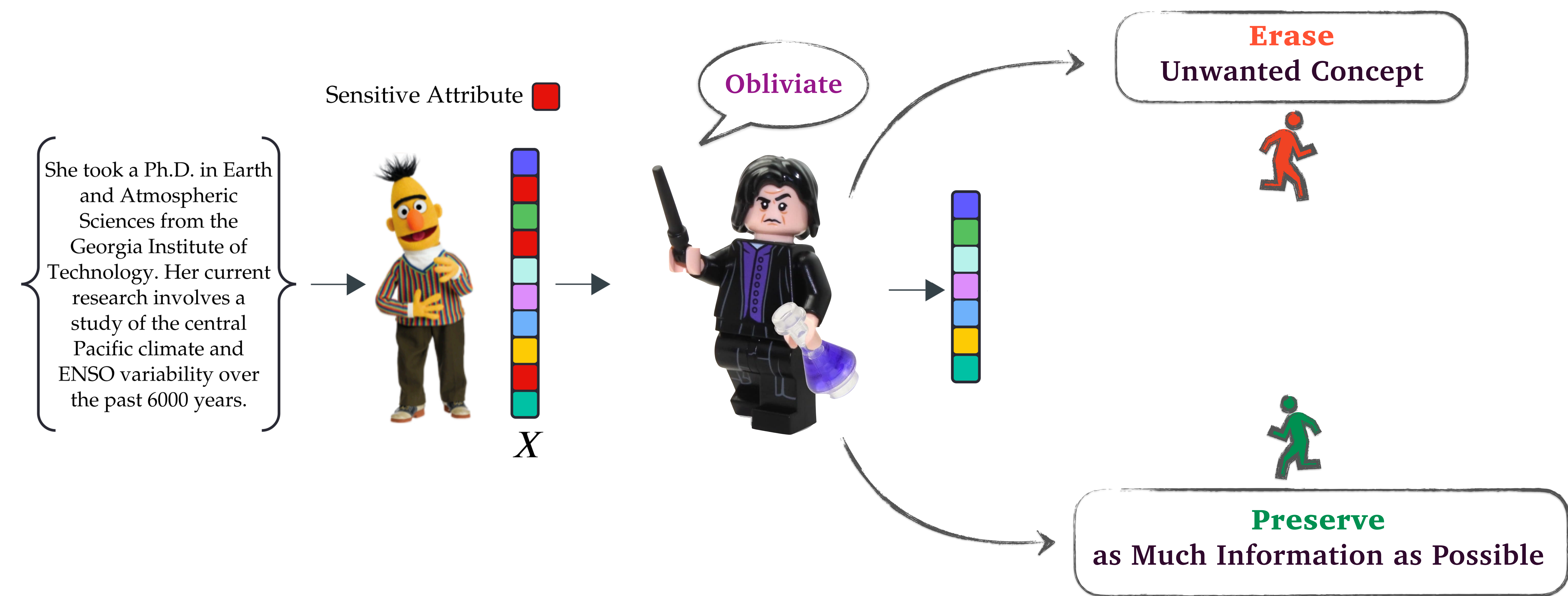
What Is The **Ultimate** Goal of Concept Erasure ?



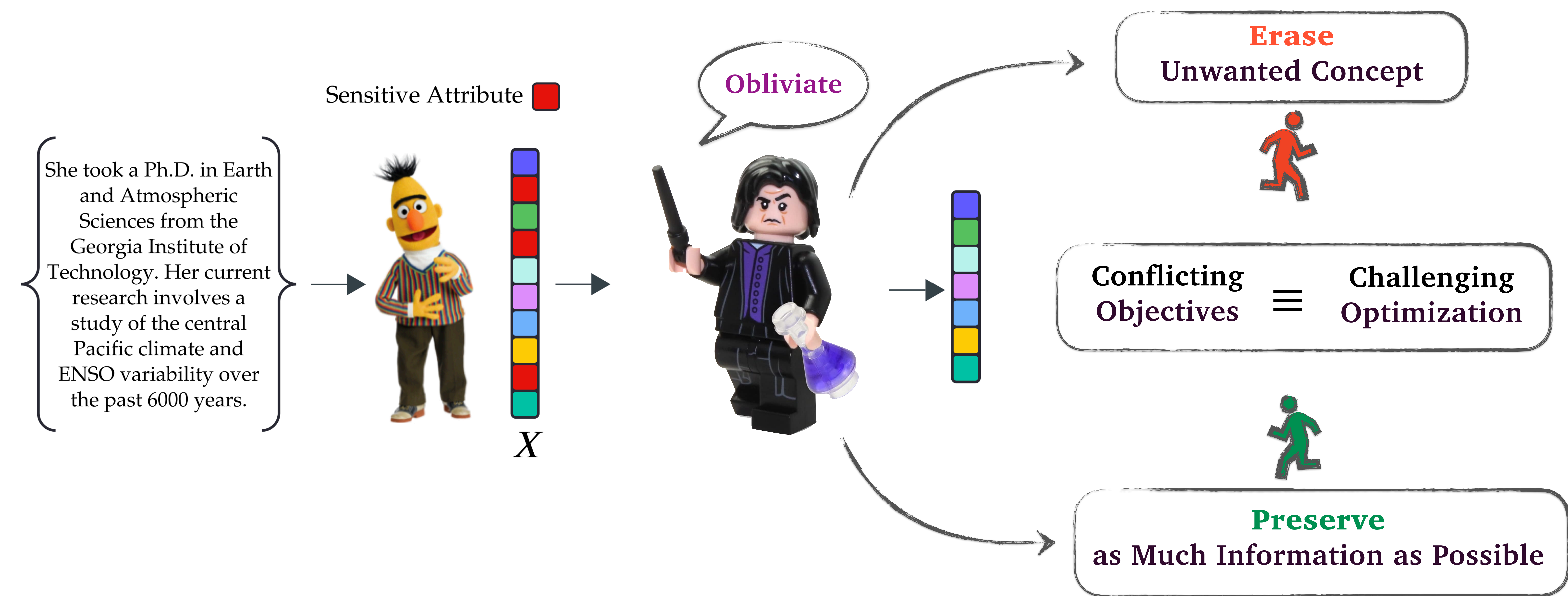
What Is The **Ultimate** Goal of Concept Erasure ?



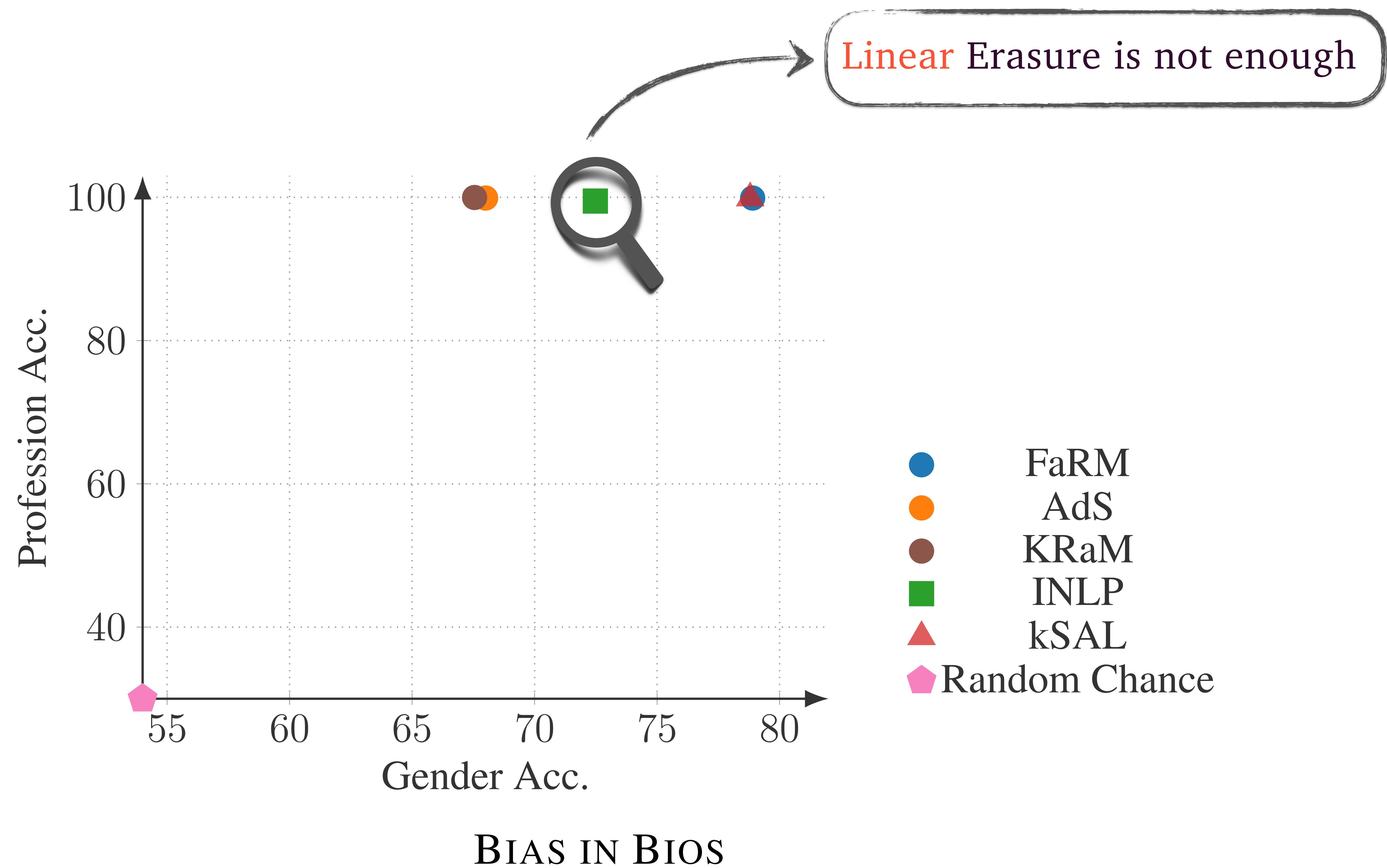
What Is The **Ultimate** Goal of Concept Erasure ?



What Is The **Ultimate** Goal of Concept Erasure ?



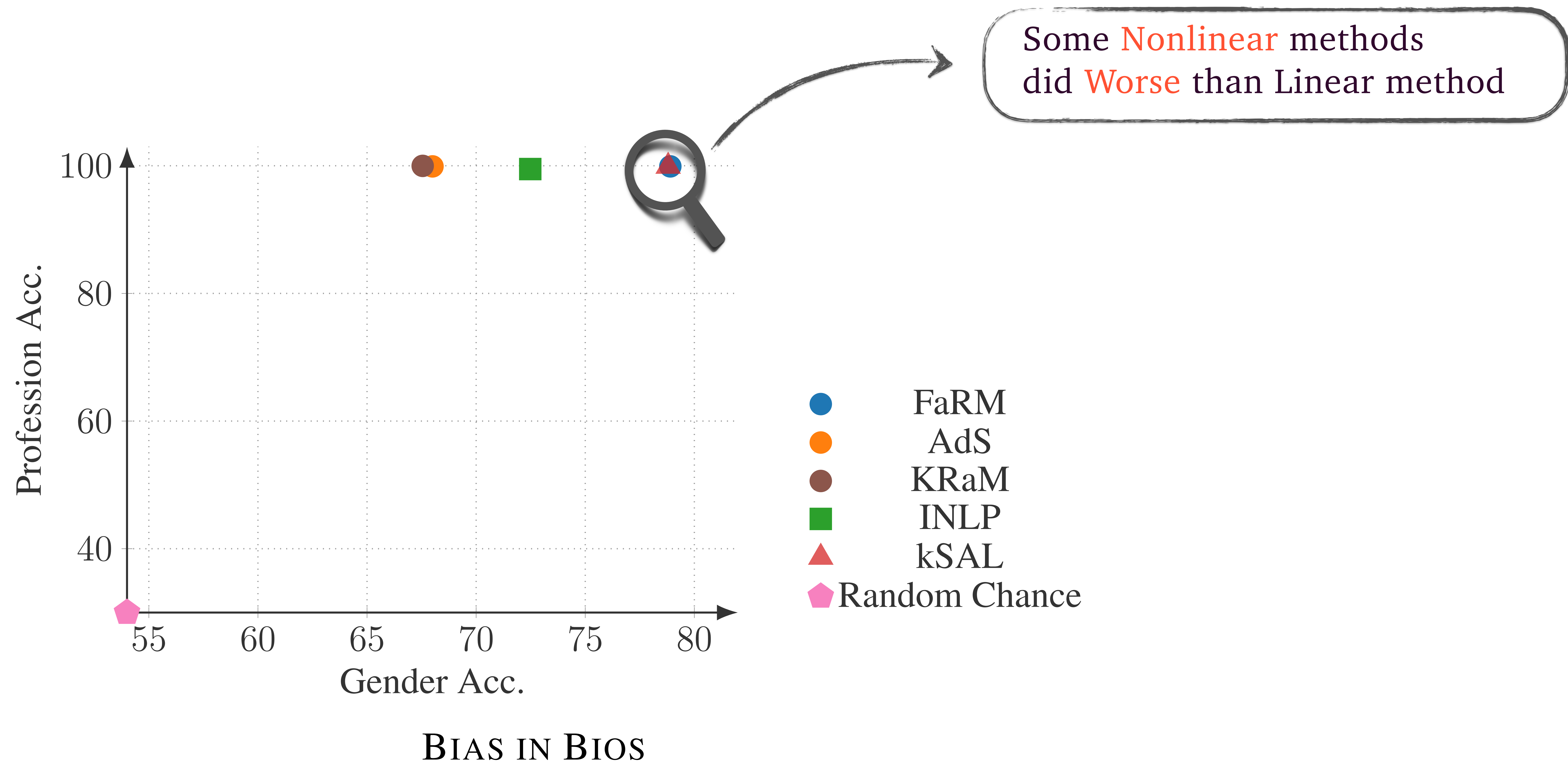
How **Effective** Are **Current** Methods In Concept Erasure ?



Utility : Profession

Unwanted : Gender

How **Effective** Are **Current** Methods In Concept Erasure ?

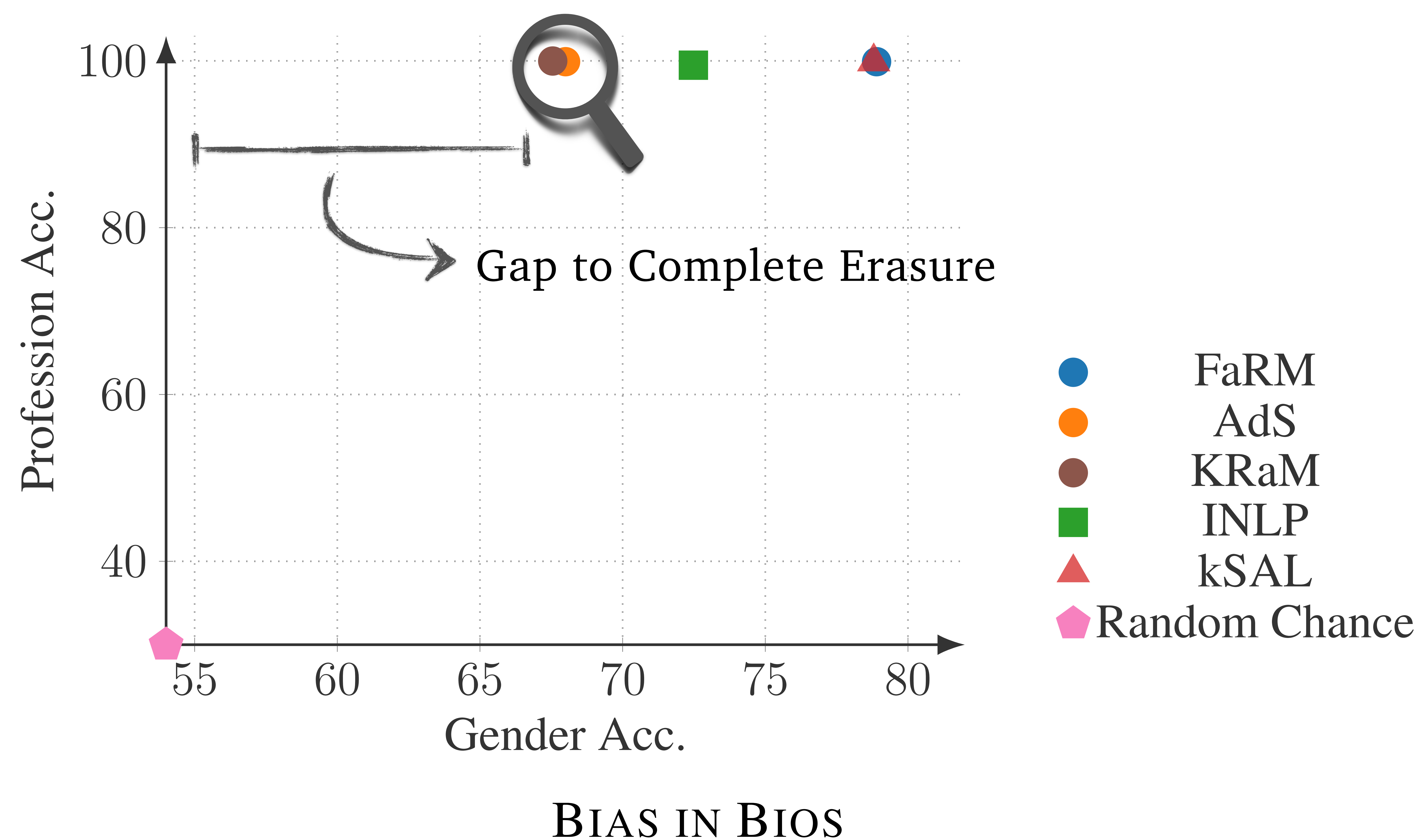


Utility : Profession

Unwanted : Gender

How **Effective** Are **Current** Methods In Concept Erasure ?

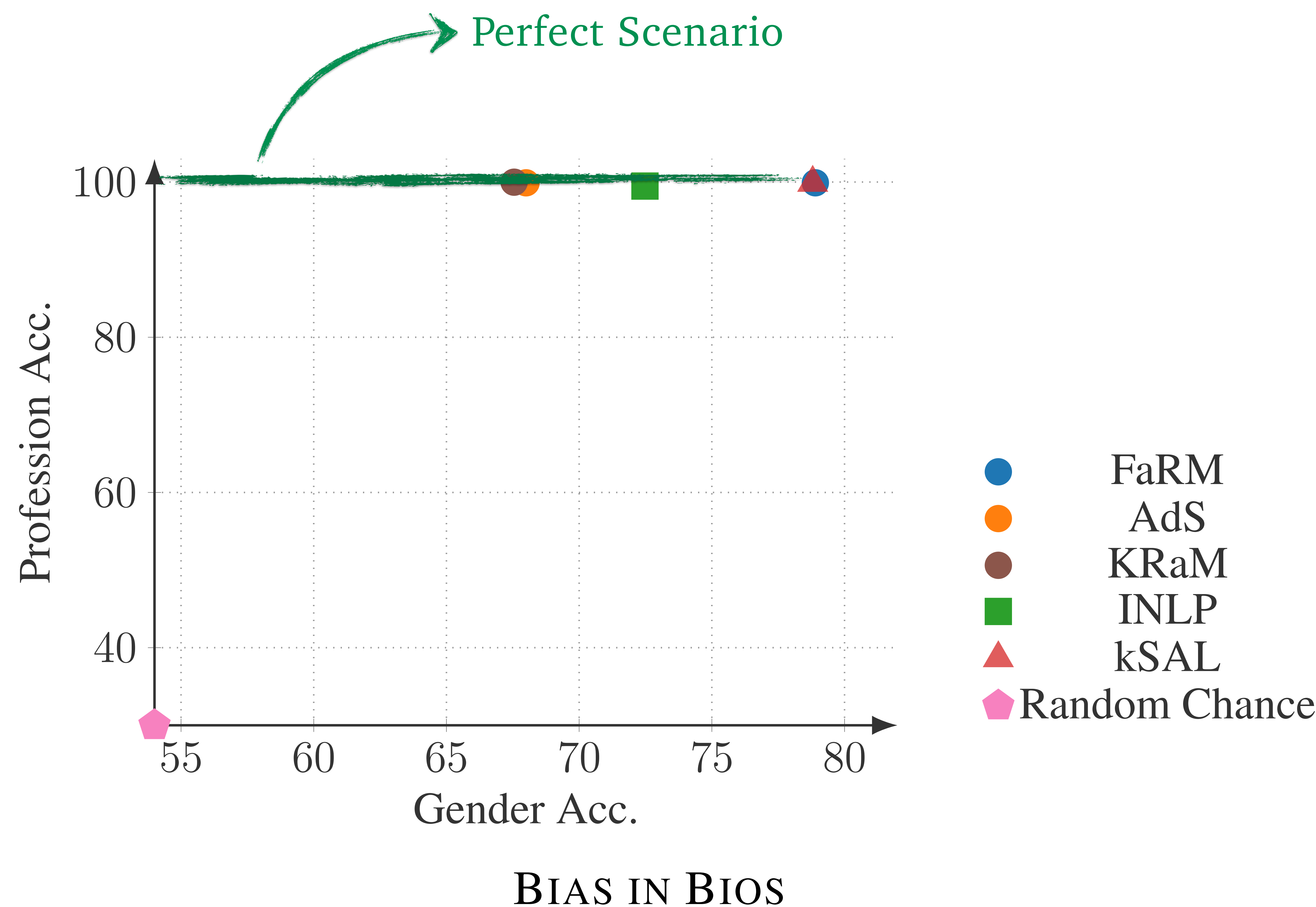
Better Performance but Erasure is **Not Complete**



Utility : Profession

Unwanted : Gender

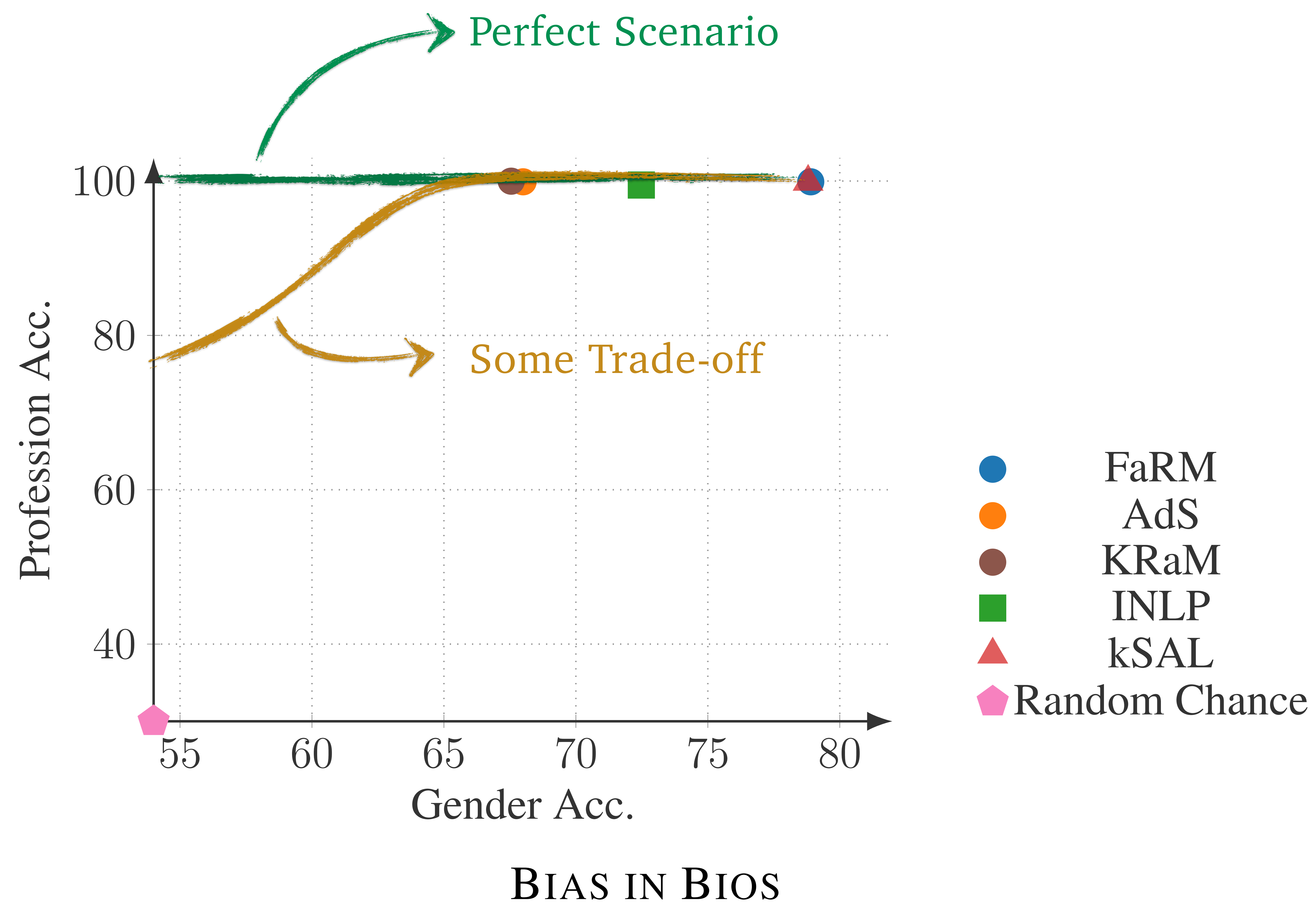
What Is The Complete **Utility-Erasure** Trade-off Profile ?



Utility : Profession

Unwanted : Gender

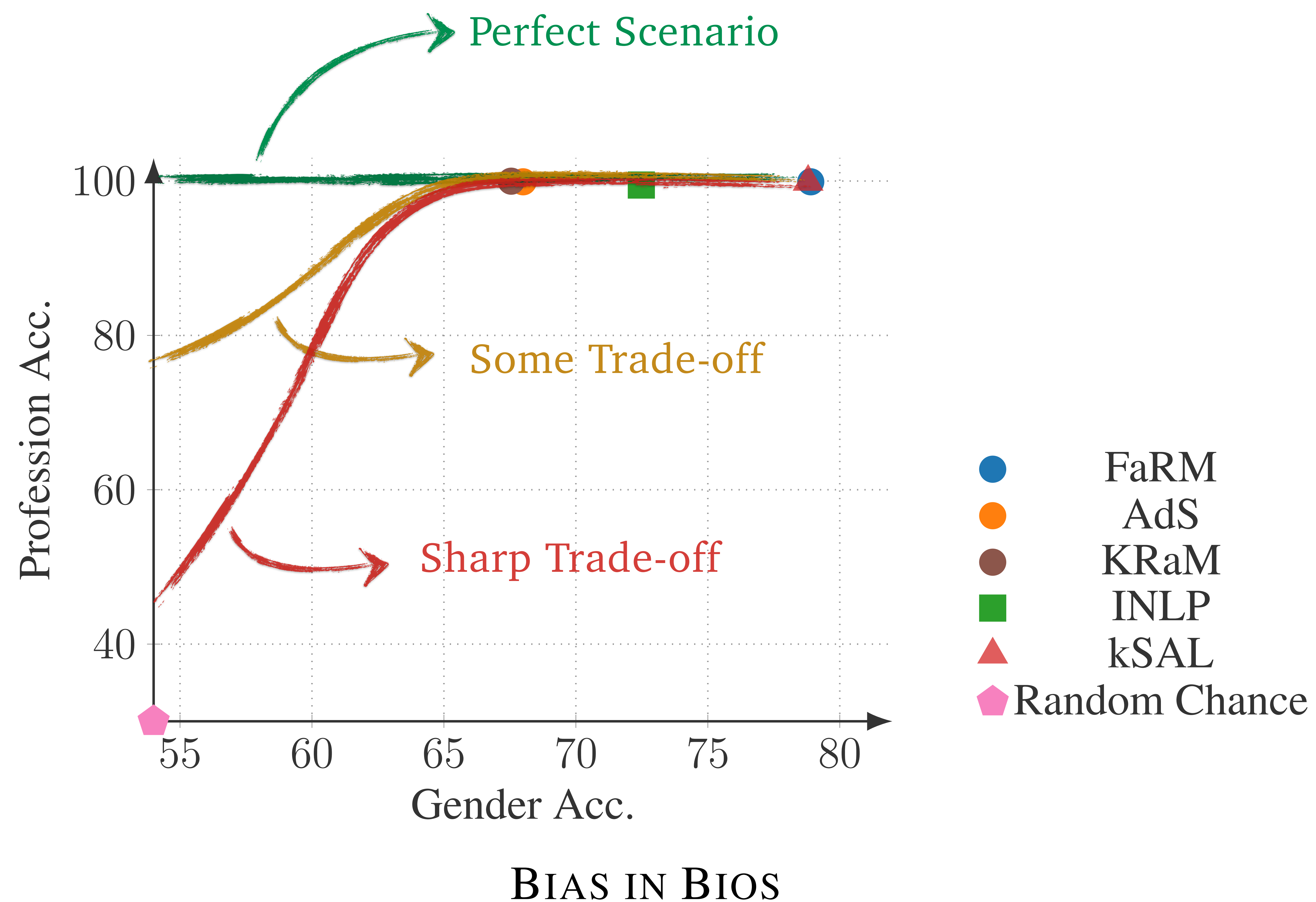
What Is The Complete **Utility-Erasure** Trade-off Profile ?



Utility : Profession

Unwanted : Gender

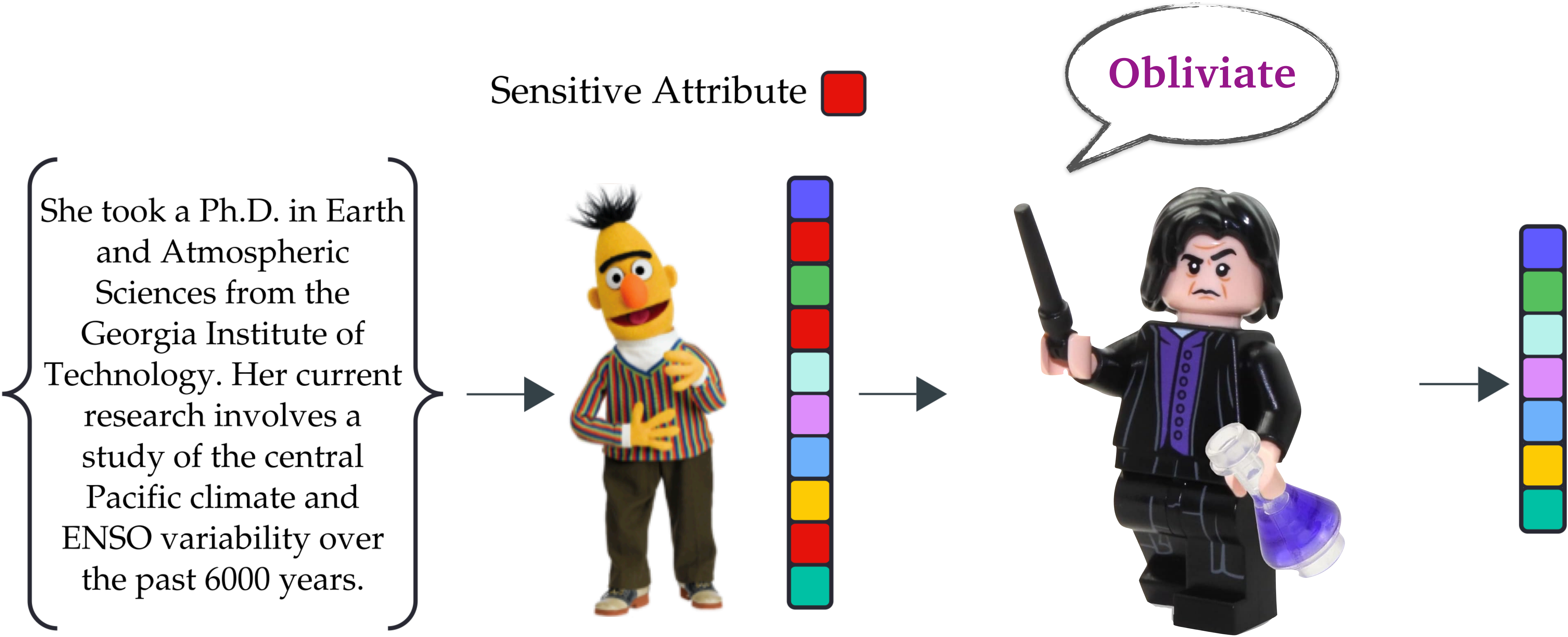
What Is The Complete **Utility-Erasure** Trade-off Profile ?



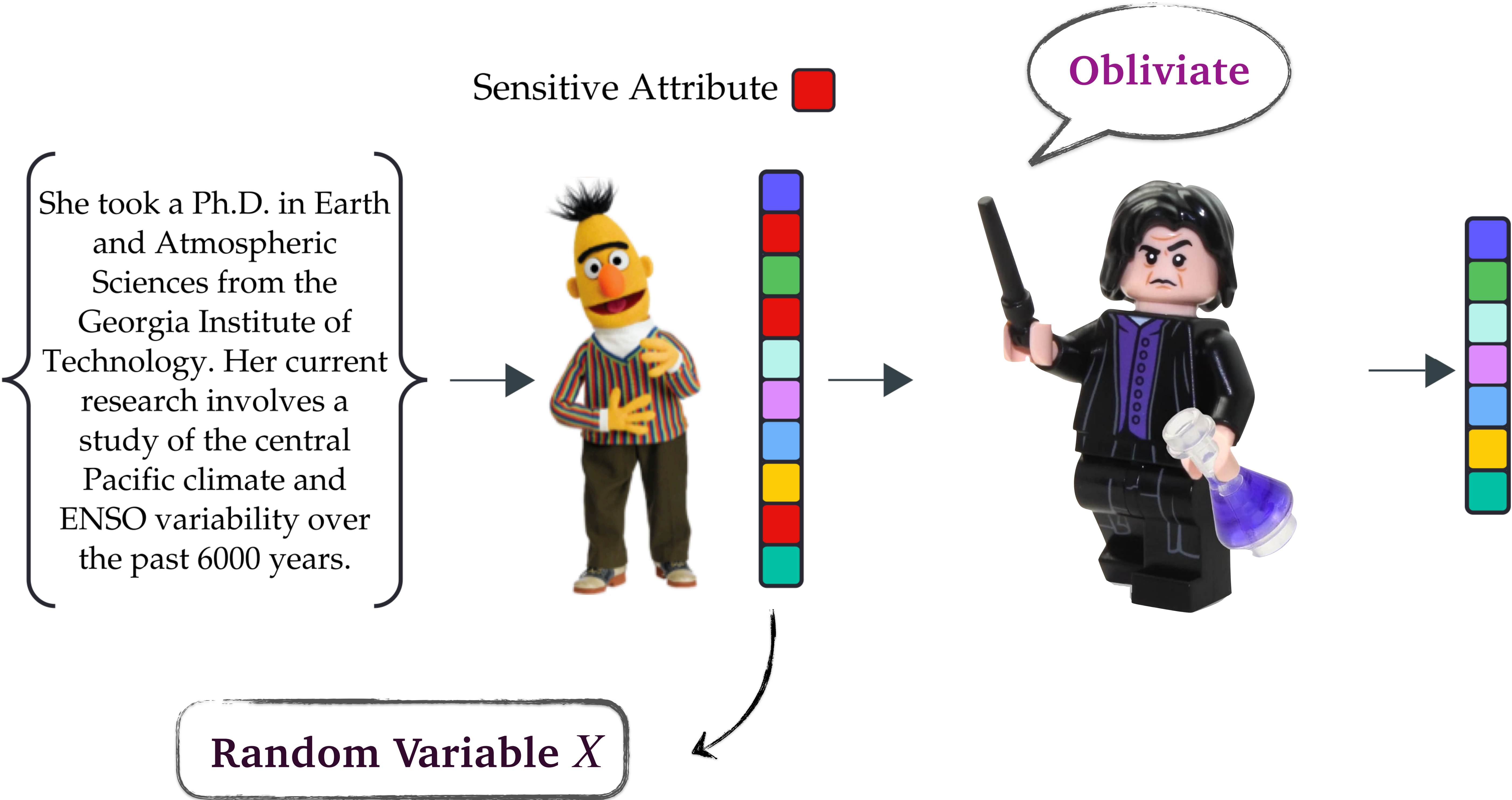
Utility : Profession

Unwanted : Gender

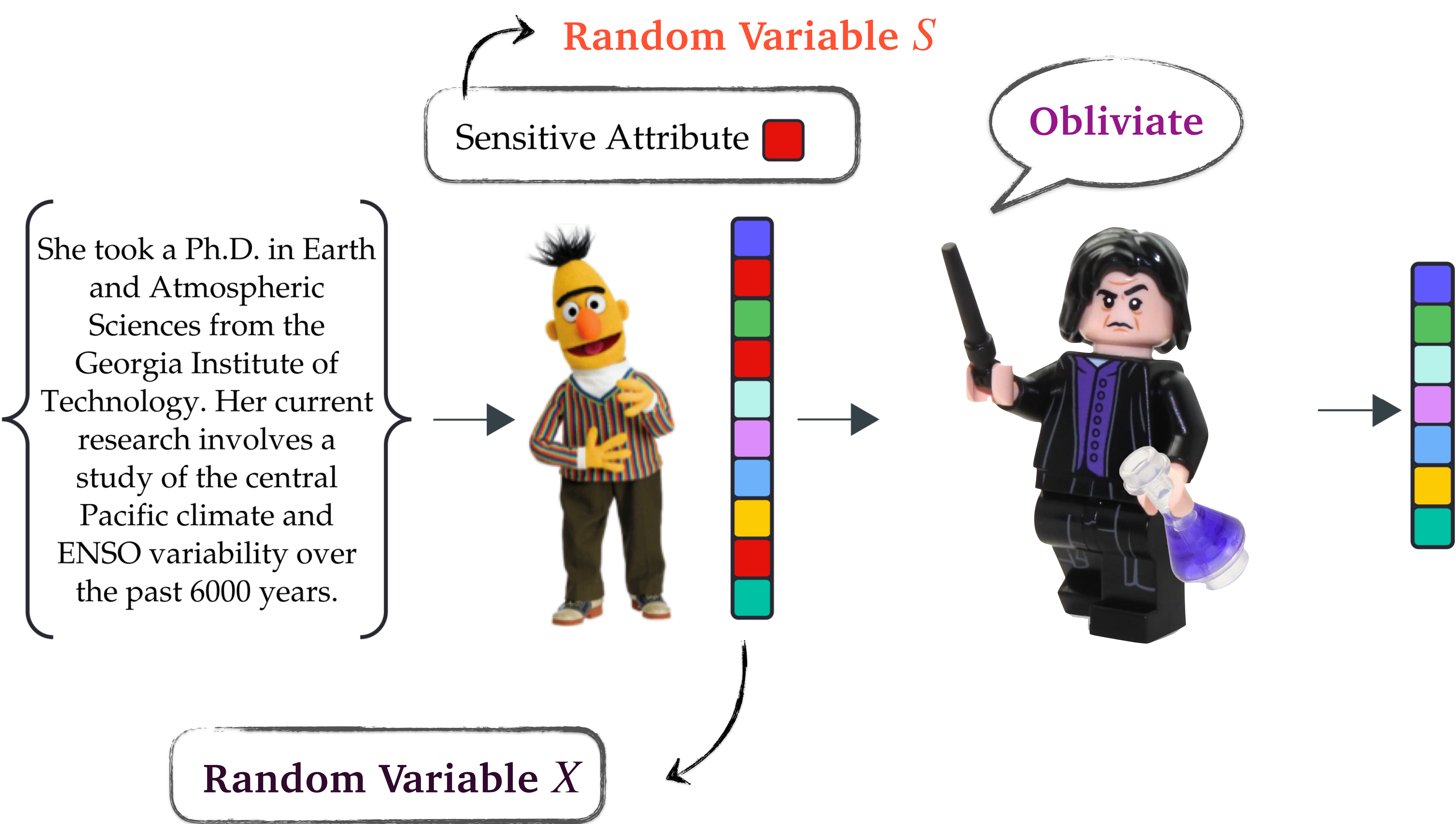
How Do **We** Formulate Concept Erasure ?



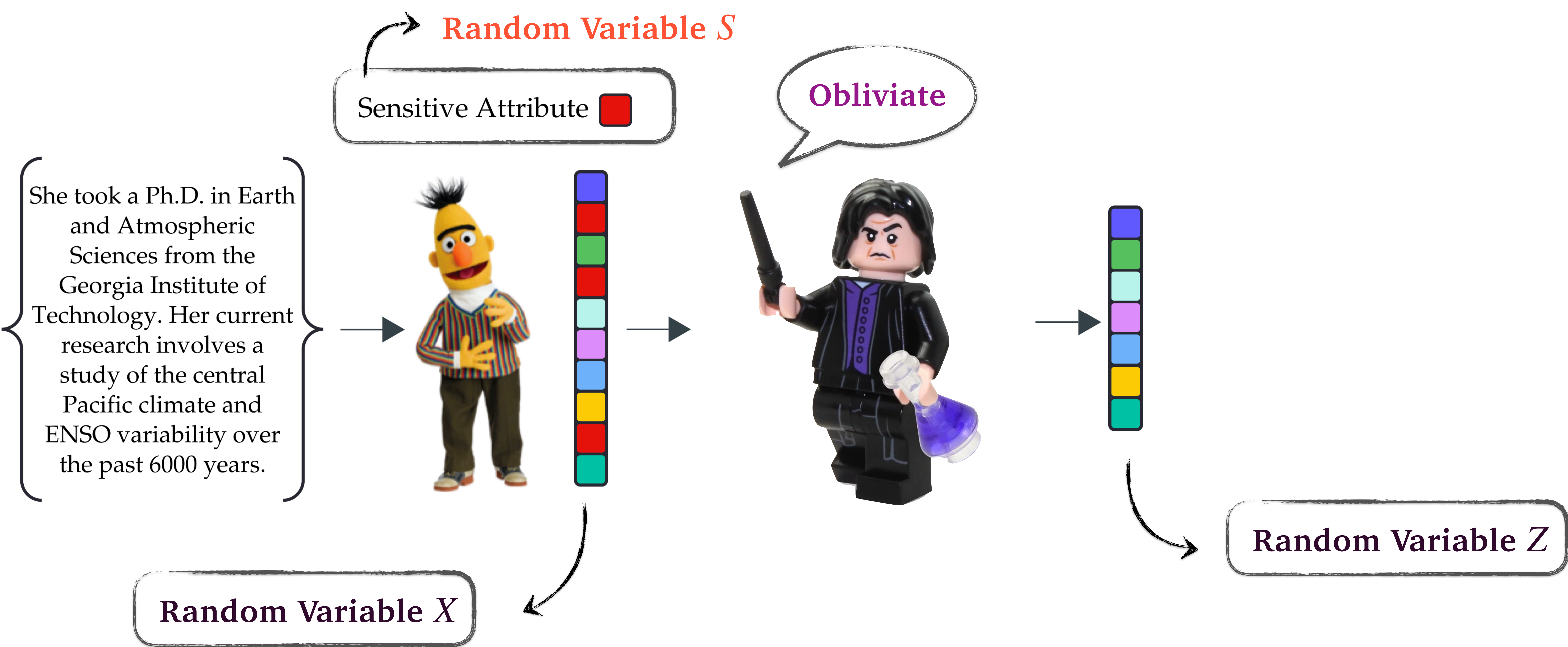
How Do **We** Formulate Concept Erasure ?



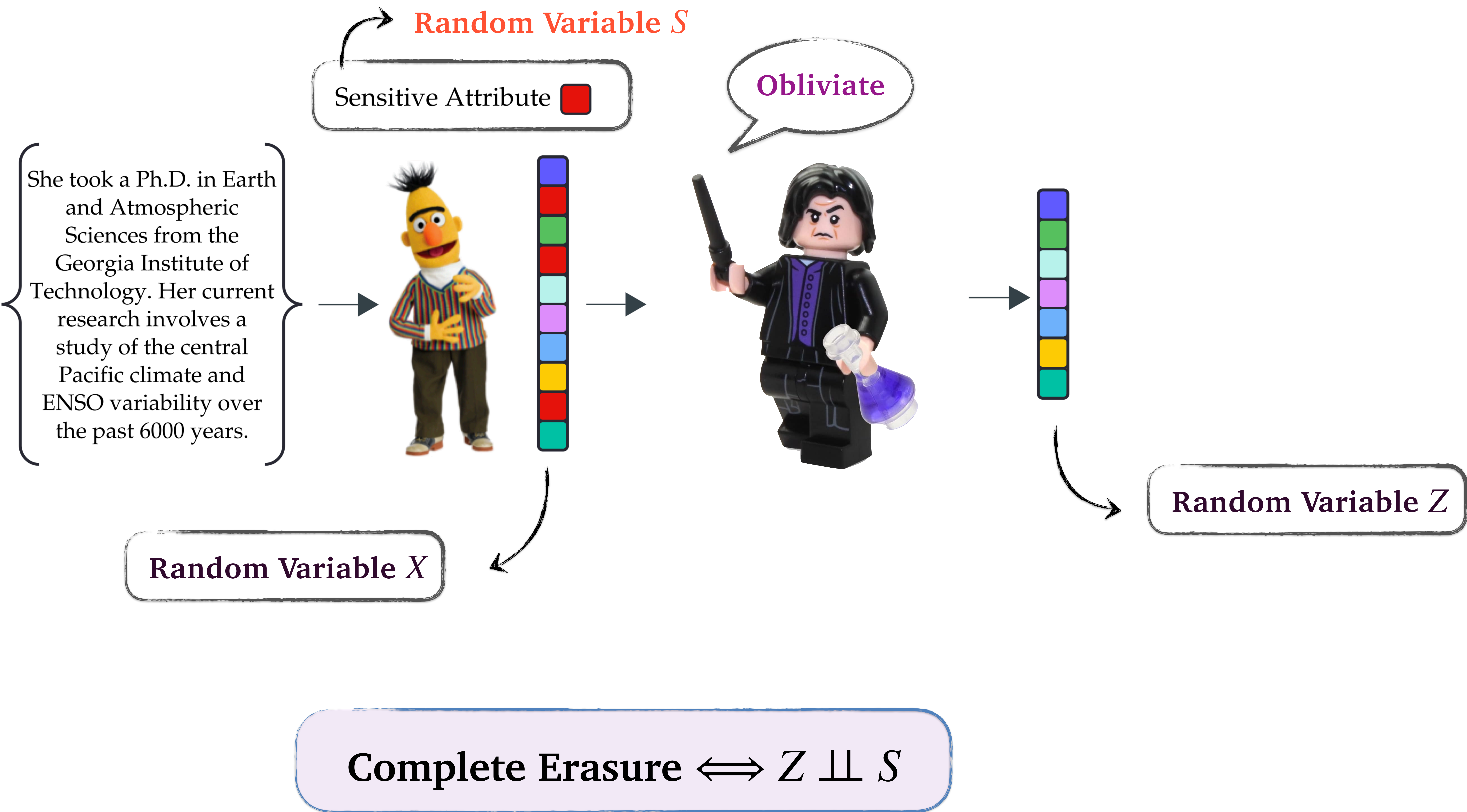
How Do **We** Formulate Concept Erasure ?



How Do **We** Formulate Concept Erasure ?

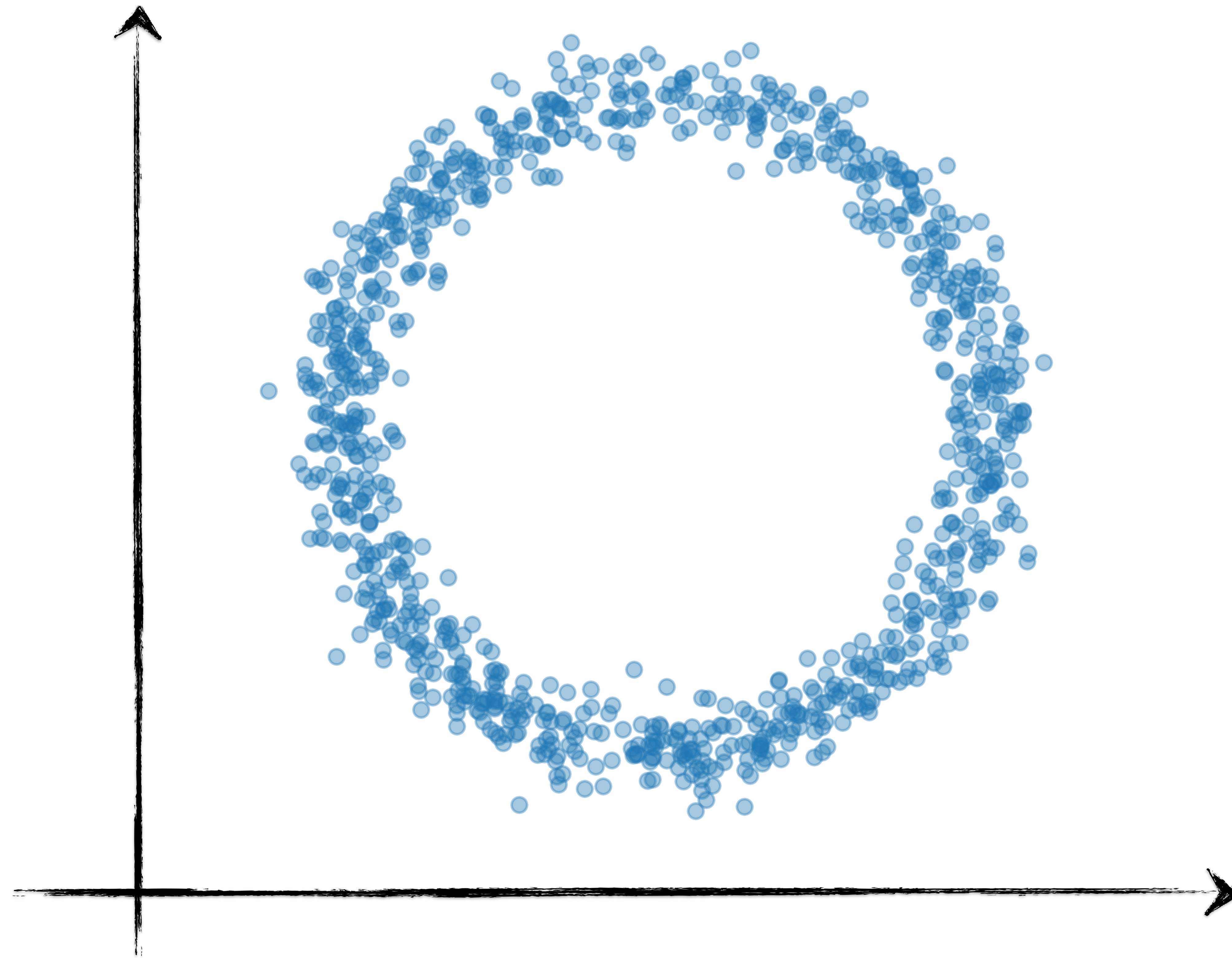


How Do **We** Formulate Concept Erasure ?



Why **Linear** Statistical Dependence Is **Not Enough** ?

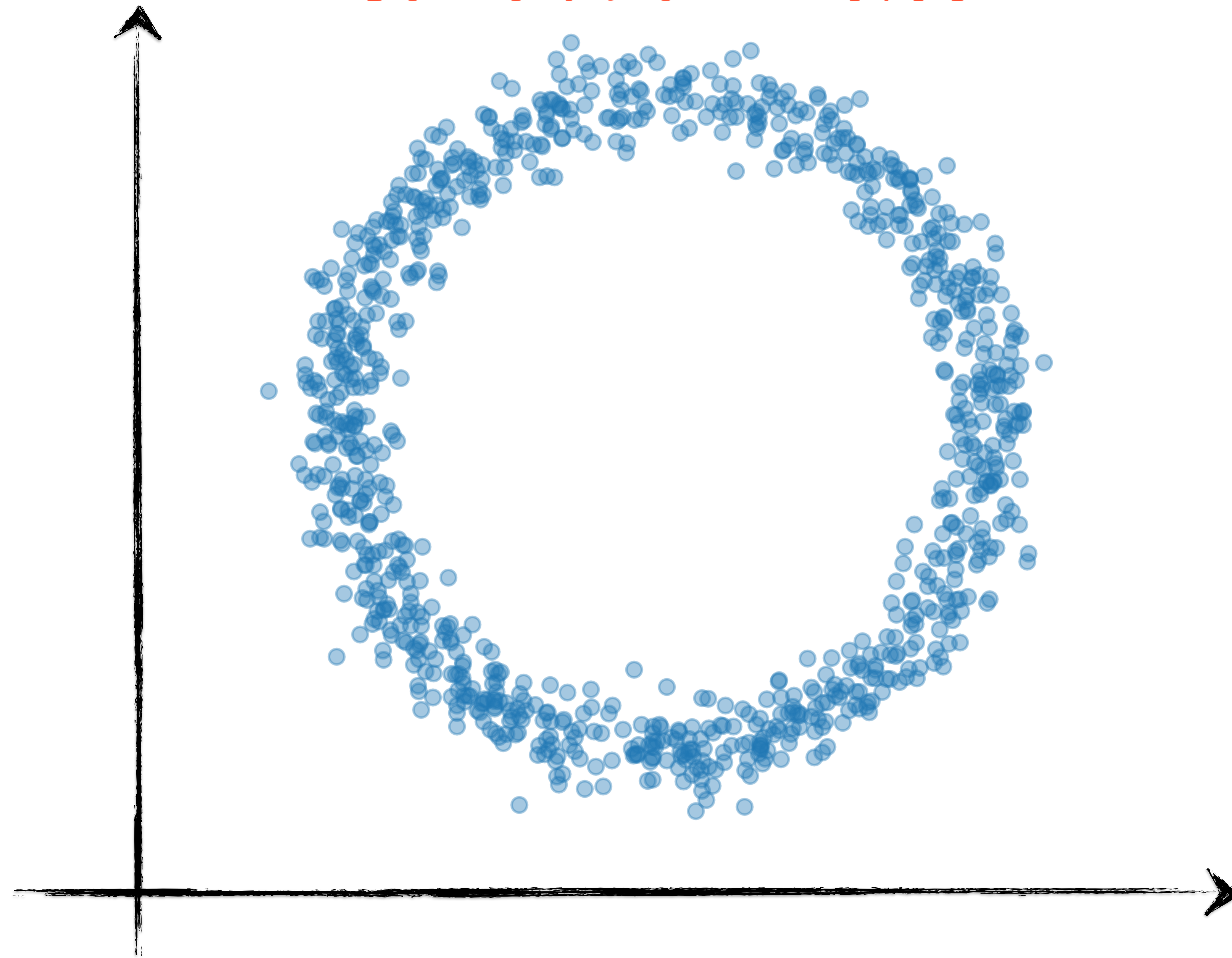
$$\text{correlation} = \frac{\text{Cov}(X, S)}{\sigma_X \sigma_S} = \frac{\mathbb{E} [(X - \mu_X)(S - \mu_S)]}{\sigma_X \sigma_S}$$



Why **Linear** Statistical Dependence Is **Not Enough** ?

$$\text{correlation} = \frac{\text{Cov}(X, S)}{\sigma_X \sigma_S} = \frac{\mathbb{E} [(X - \mu_X)(S - \mu_S)]}{\sigma_X \sigma_S}$$

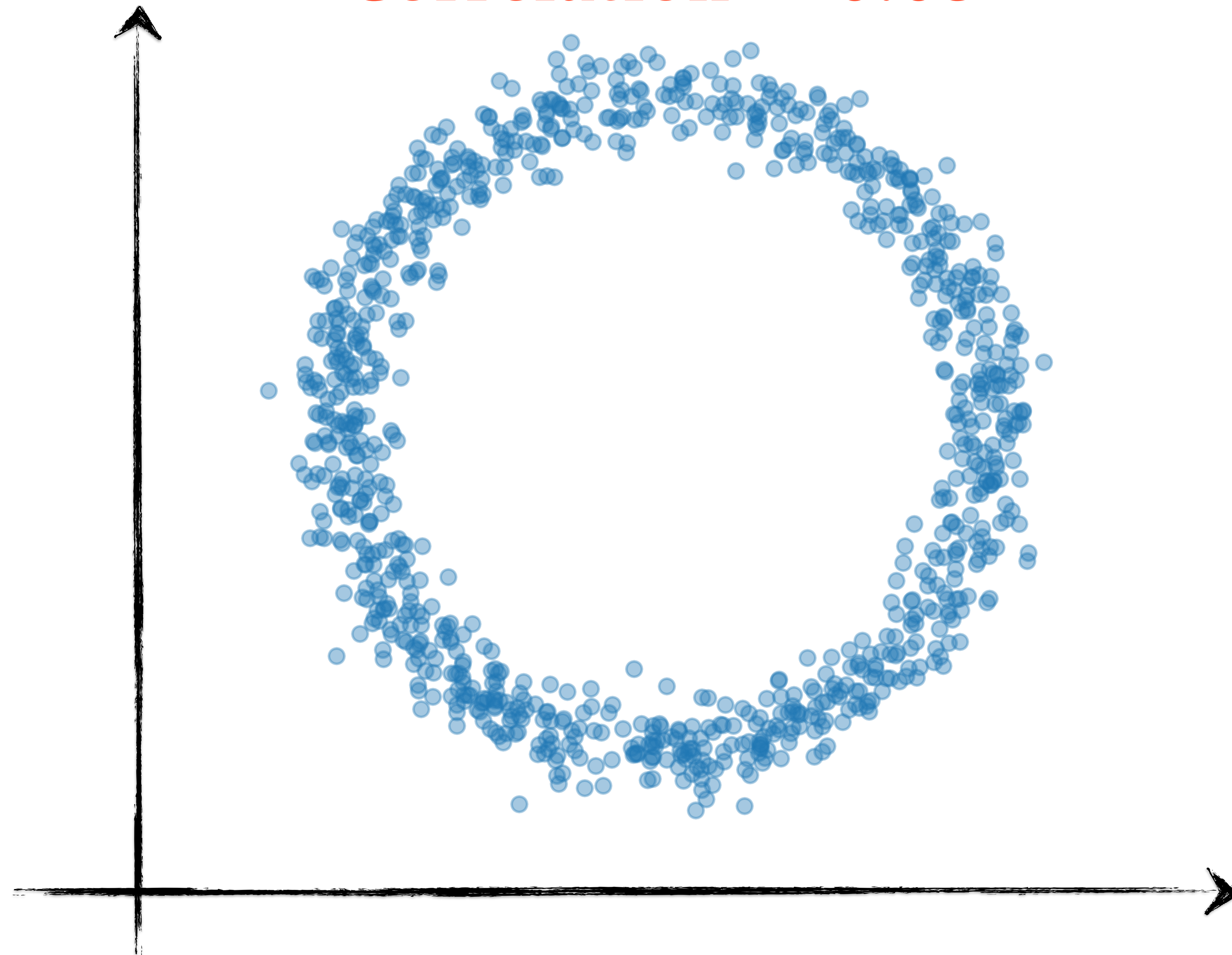
Correlation = 0.03



Why **Linear** Statistical Dependence Is **Not Enough** ?

$$\text{correlation} = \frac{\text{Cov}(X, S)}{\sigma_X \sigma_S} = \frac{\mathbb{E} [(X - \mu_X)(S - \mu_S)]}{\sigma_X \sigma_S}$$

Correlation = 0.03

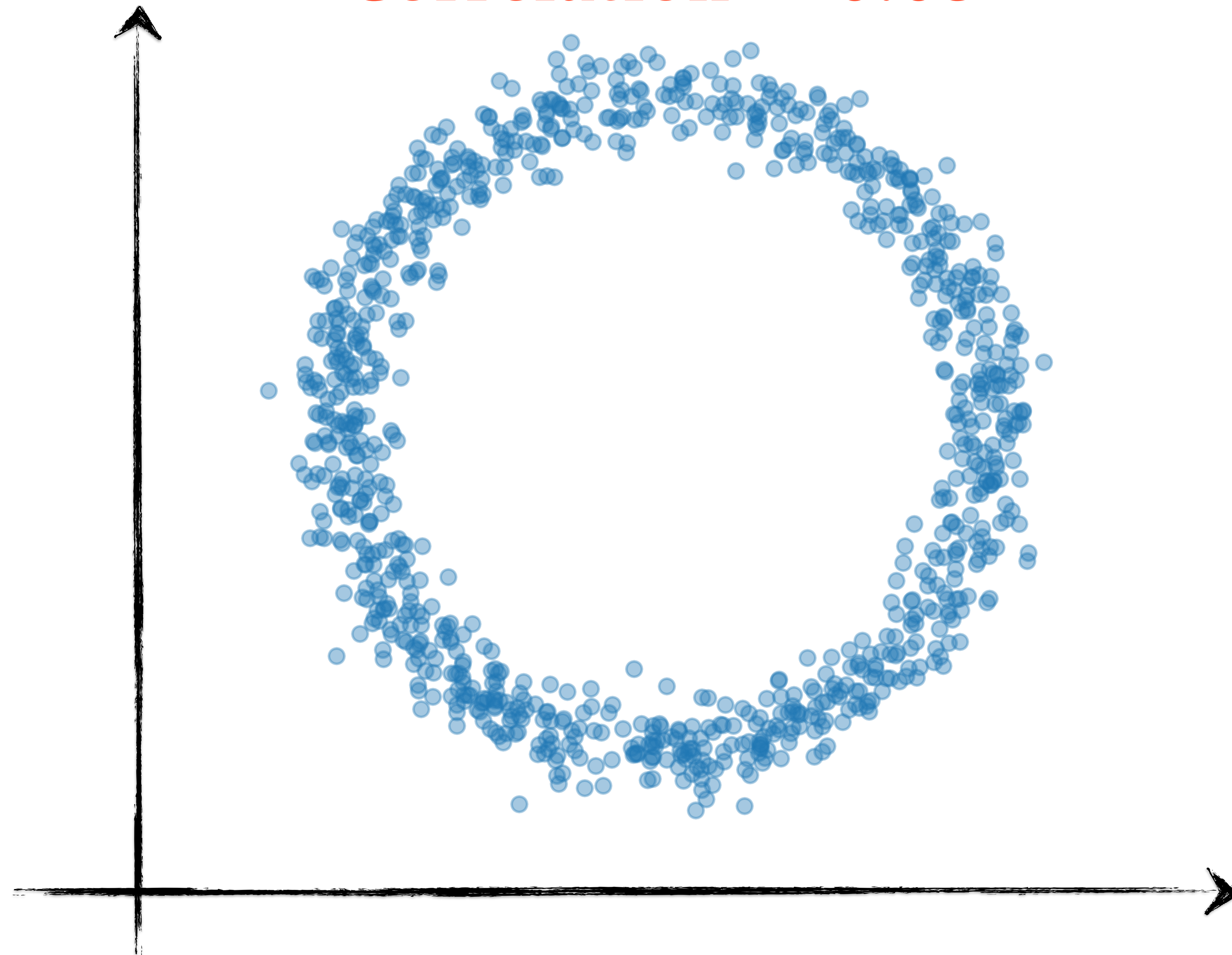


Independence \rightarrow Correlation = 0

Why **Linear** Statistical Dependence Is **Not Enough** ?

$$\text{correlation} = \frac{\text{Cov}(X, S)}{\sigma_X \sigma_S} = \frac{\mathbb{E} [(X - \mu_X)(S - \mu_S)]}{\sigma_X \sigma_S}$$

Correlation = 0.03

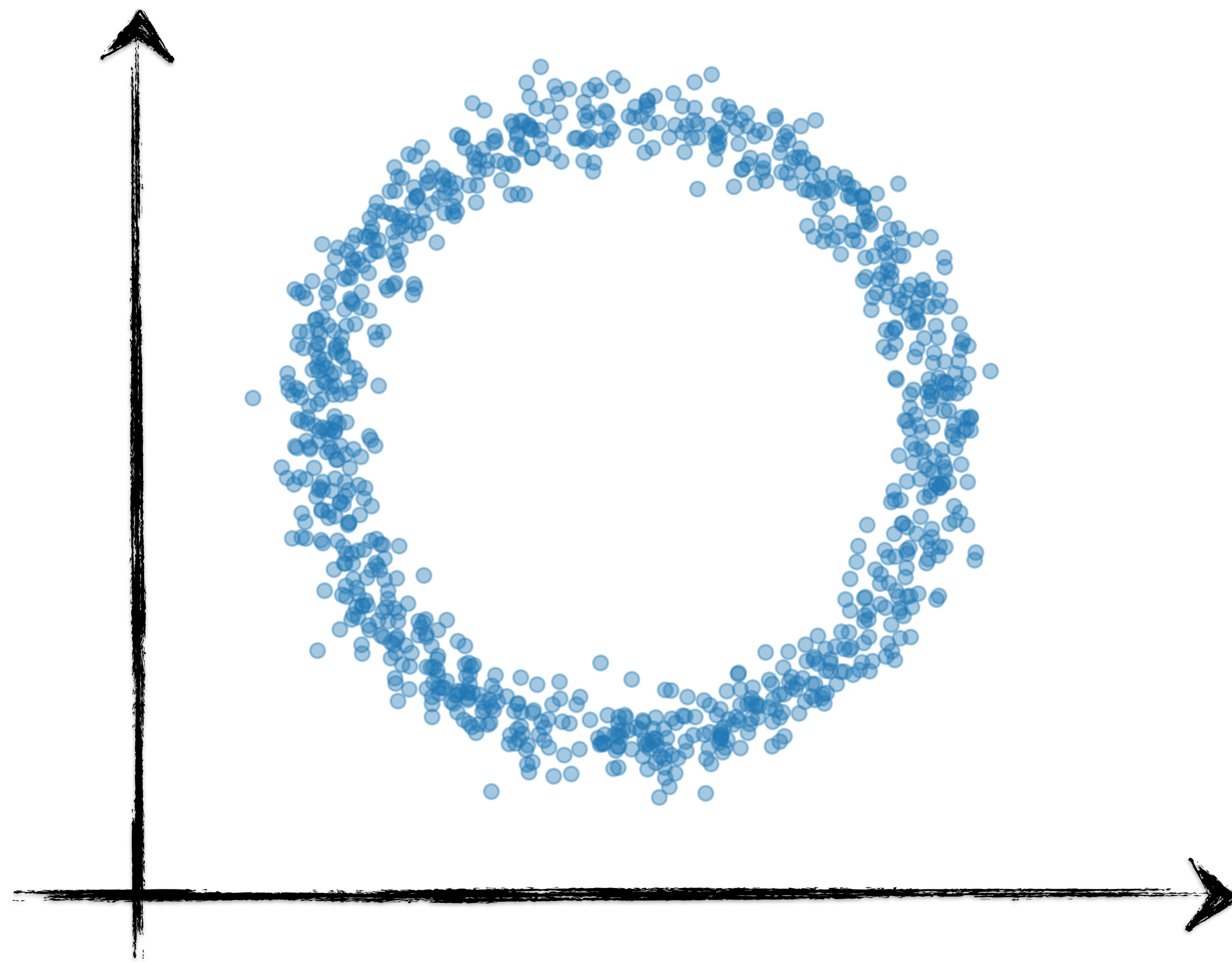


Independence \rightarrow Correlation = 0

Correlation = 0 \nrightarrow Independence

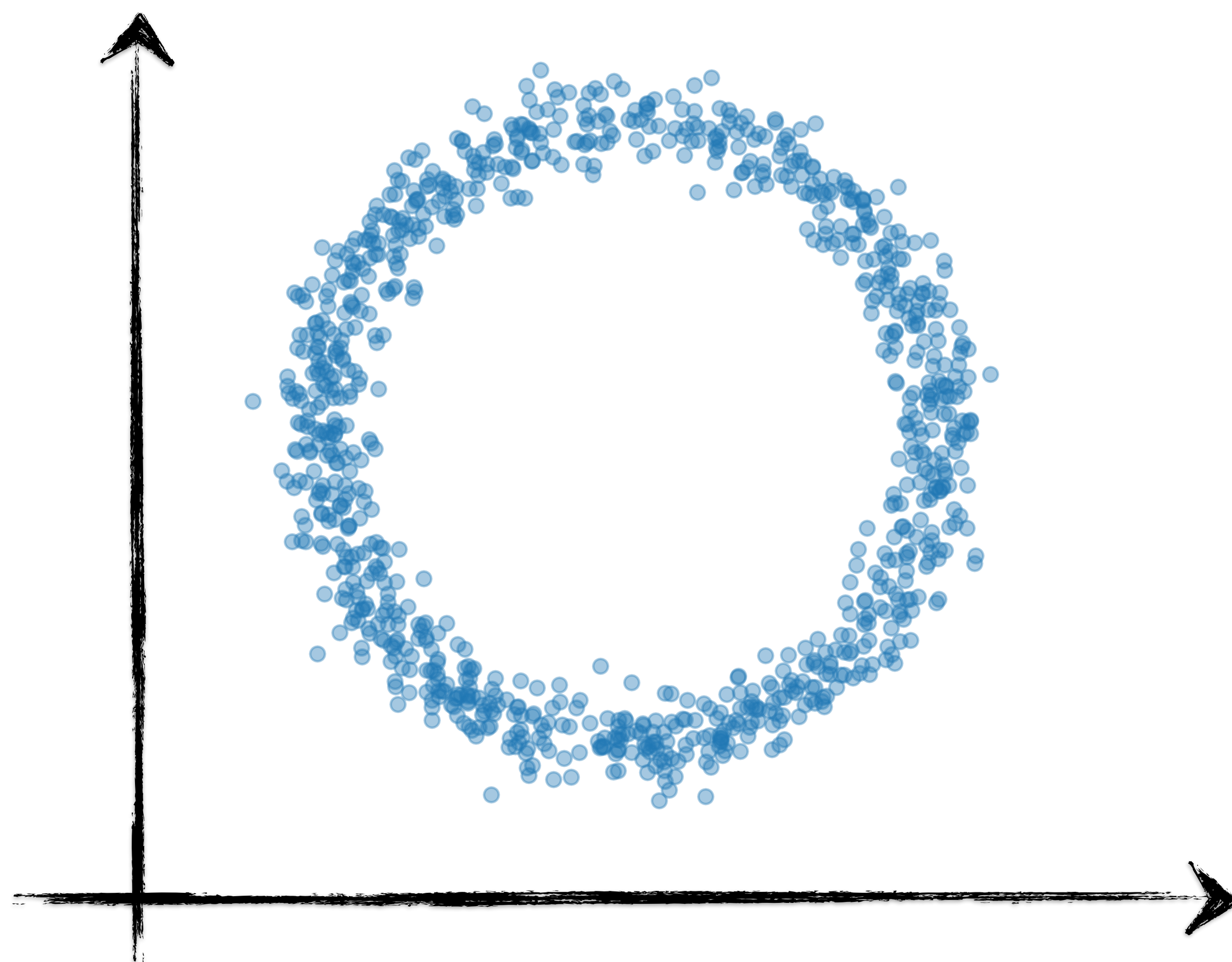
How Do **We** Capture **Nonlinear** Statistical Dependency ?

Correlation = 0.03



How Do **We** Capture **Nonlinear** Statistical Dependency ?

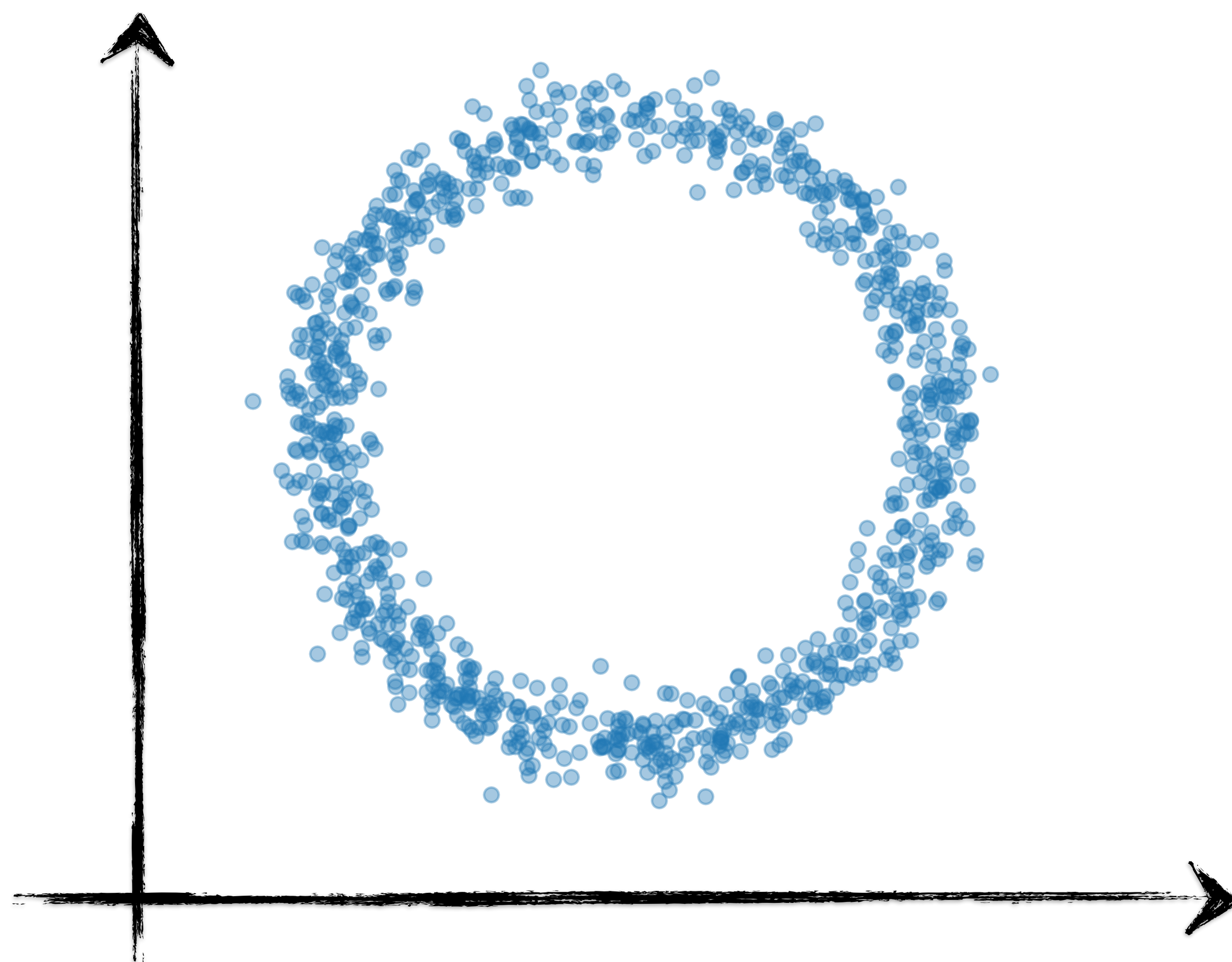
Correlation = 0.03



Transform Nonlinear Dependency between X and S

How Do **We** Capture **Nonlinear** Statistical Dependency ?

Correlation = 0.03

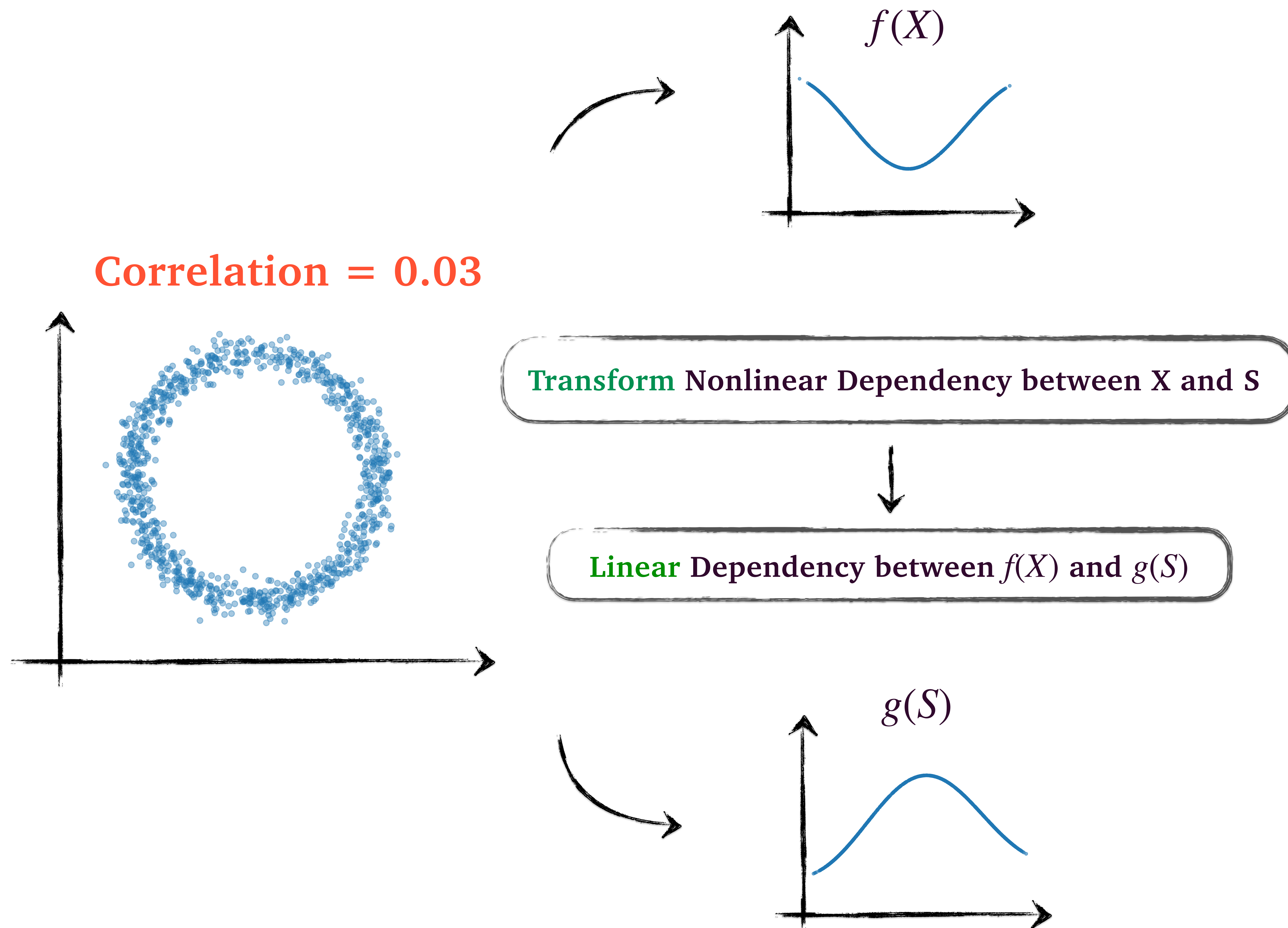


Transform Nonlinear Dependency between X and S

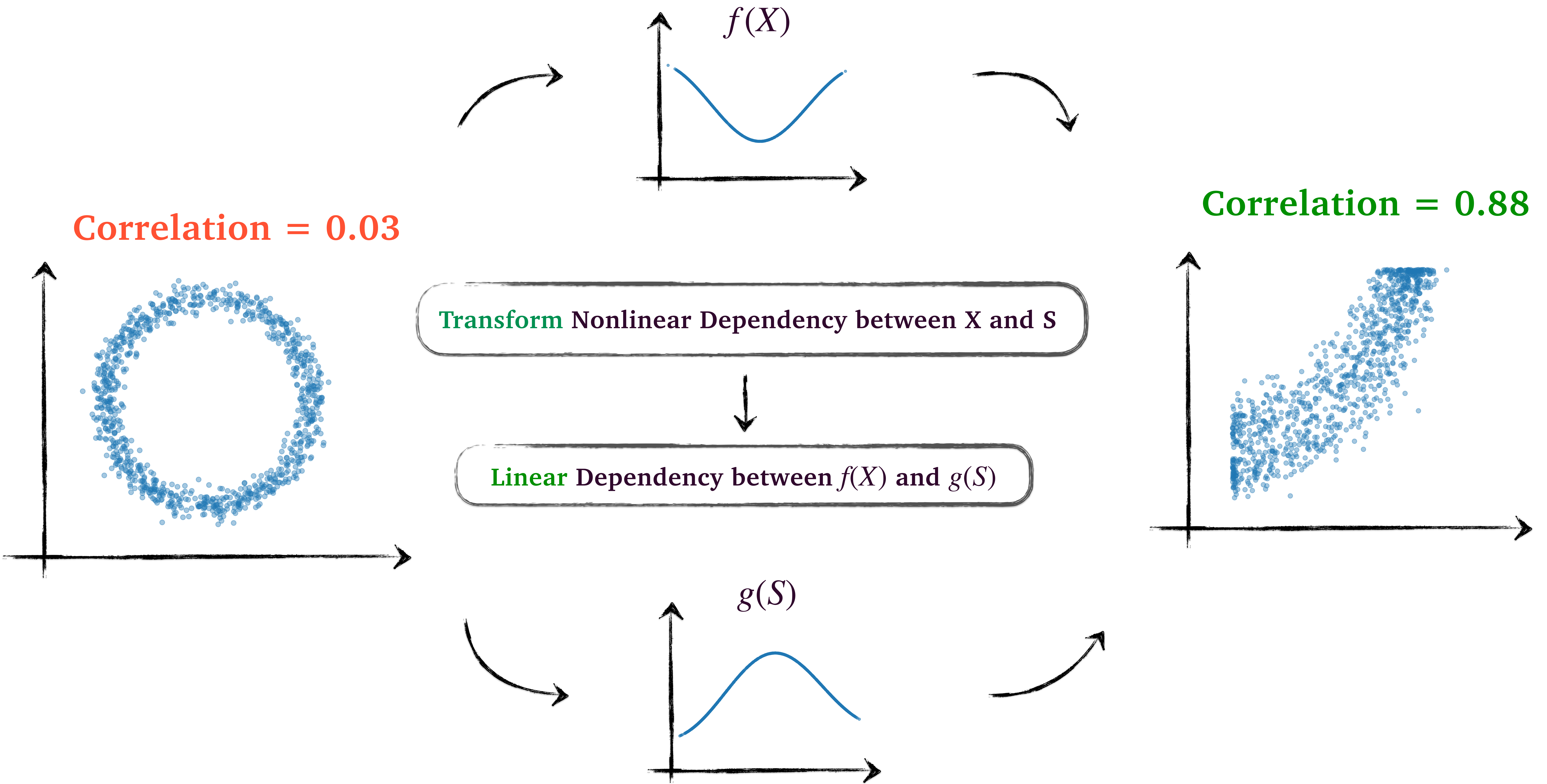


Linear Dependency between $f(X)$ and $g(S)$

How Do **We** Capture **Nonlinear** Statistical Dependency ?



How Do **We** Capture **Nonlinear** Statistical Dependency ?



How Do We **Formalize** a **Nonlinear** Statistical Dependency Test ?

Let \mathcal{F} and \mathcal{G} be a characteristic RKHS

How Do We **Formalize** a **Nonlinear** Statistical Dependency Test ?

Let \mathcal{F} and \mathcal{G} be a characteristic RKHS

$$\sup_{g \in \mathcal{G}} \sup_{f \in \mathcal{F}} \mathbf{Cov}(f(X), g(S)) = \langle g, \mathbf{Cov} f \rangle_{\mathcal{G}} \quad s.t. \quad \|f\|_{\mathcal{F}} = \|g\|_{\mathcal{G}} = 1$$

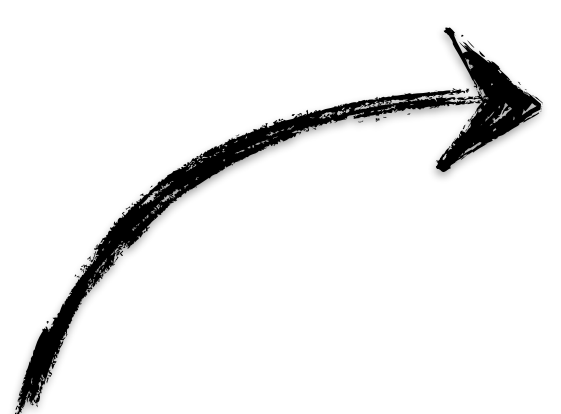
Regularization
Constraint



How Do We **Formalize** a **Nonlinear** Statistical Dependency Test ?

Let \mathcal{F} and \mathcal{G} be a characteristic RKHS

$$\sup_{g \in \mathcal{G}} \sup_{f \in \mathcal{F}} \mathbf{Cov}(f(X), g(S)) = \langle g, \mathbb{Cov} f \rangle_{\mathcal{G}} \quad s.t. \quad \|f\|_{\mathcal{F}} = \|g\|_{\mathcal{G}} = 1$$

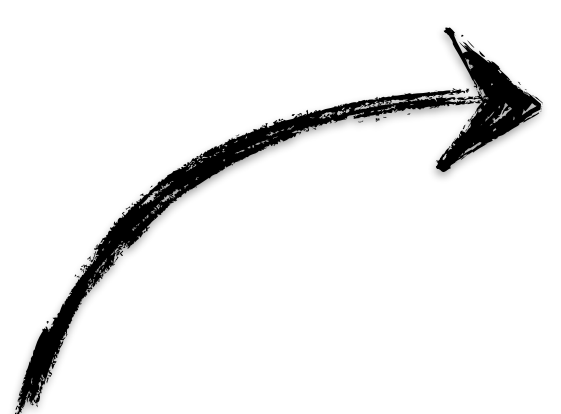
Regularization
Constraint 

$\|\mathbb{Cov}\|_{HS}^2$ **Hilbert-Schmidt Norm** of Covariance Operator
Indicates **Existence** of Such Functions

How Do We **Formalize** a **Nonlinear** Statistical Dependency Test ?

Let \mathcal{F} and \mathcal{G} be a characteristic RKHS

$$\sup_{g \in \mathcal{G}} \sup_{f \in \mathcal{F}} \mathbf{Cov}(f(X), g(S)) = \langle g, \mathbb{Cov} f \rangle_{\mathcal{G}} \quad s.t. \quad \|f\|_{\mathcal{F}} = \|g\|_{\mathcal{G}} = 1$$

Regularization
Constraint 

$\|\mathbb{Cov}\|_{HS}^2$ **Hilbert-Schmidt Norm** of Covariance Operator
Indicates **Existence** of Such Functions

$$\|\mathbb{Cov}\|_{HS}^2 = 0 \iff \text{Statistical Independence}$$

How Do We **Formalize** a **Nonlinear** Statistical Dependency Test ?

Let \mathcal{F} and \mathcal{G} be a characteristic RKHS

$$\sup_{g \in \mathcal{G}} \sup_{f \in \mathcal{F}} \mathbf{Cov}(f(X), g(S)) = \langle g, \mathbb{Cov} f \rangle_{\mathcal{G}} \quad s.t. \quad \|f\|_{\mathcal{F}} = \|g\|_{\mathcal{G}} = 1$$

Regularization
Constraint

$\|\mathbb{Cov}\|_{HS}^2$ **Hilbert-Schmidt Norm** of Covariance Operator
Indicates **Existence** of Such Functions

$$\|\mathbb{Cov}\|_{HS}^2 = 0 \iff \text{Statistical Independence}$$

The Correct Proxy To Impose Statistical Independence

How Do We **Formulate** Erasure Using Statistical Independence ?

$$Z_{\theta} = \varepsilon(X; \theta) \rightarrow$$



How Do We **Formulate** Erasure Using Statistical Independence ?

$$Z_{\theta} = \varepsilon(X; \theta) \rightarrow$$



- Z_{θ} \rightarrow Random Variable After Erasure
- S \rightarrow Random Variable Represent **Unwanted** Label
- Y \rightarrow Random Variable Represent **Utility** Label

How Do We **Formulate** Erasure Using Statistical Independence ?

$$Z_{\theta} = \varepsilon(X; \theta) \rightarrow$$



Z_{θ} \rightarrow Random Variable After Erasure

S \rightarrow Random Variable Represent **Unwanted** Label

Y \rightarrow Random Variable Represent **Utility** Label

$$\inf_{\theta} \text{HSIC}(Z_{\theta}, S) - \text{HSIC}(Z_{\theta}, Y)$$

Minimize Statistical Dependency

A Proxy to Preserve Utility Information

How Do We **Formulate** Erasure Using Statistical Independence ?



$$Z_{\theta} = \varepsilon(X; \theta) \rightarrow$$



Z_{θ} \rightarrow Random Variable After Erasure

S \rightarrow Random Variable Represent **Unwanted** Label

Y \rightarrow Random Variable Represent **Utility** Label


$$\inf_{\theta} \text{HSIC}(Z_{\theta}, S) - \text{HSIC}(Z_{\theta}, Y)$$

Conflicting
Objectives \equiv Challenging
Optimization

Minimize Statistical Dependency

A Proxy to Preserve Utility Information

How Do We **Solve** This Optimization ?

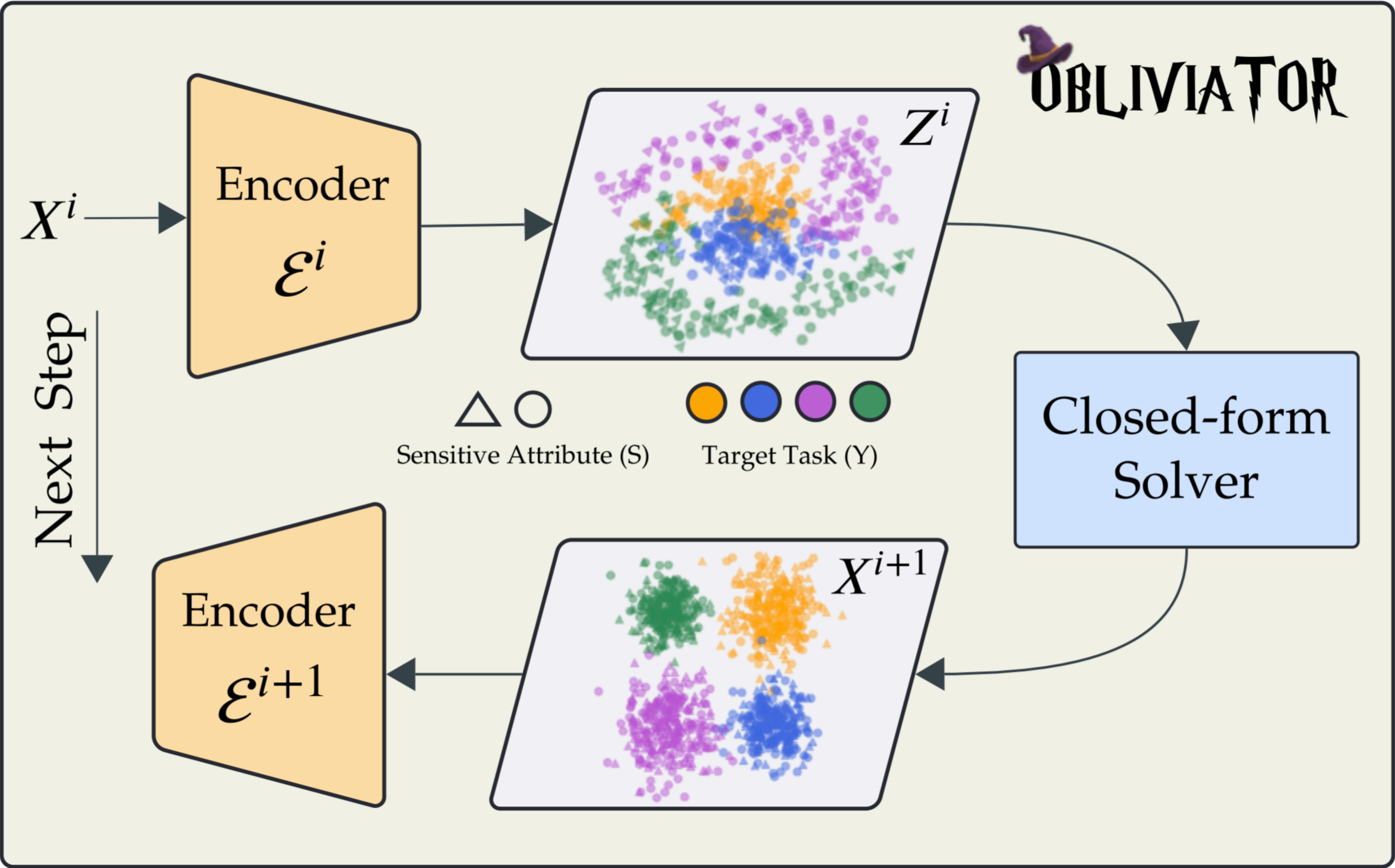
Single-Shot Optimization Leads to **Poor** Solutions

How Do We **Solve** This Optimization ?

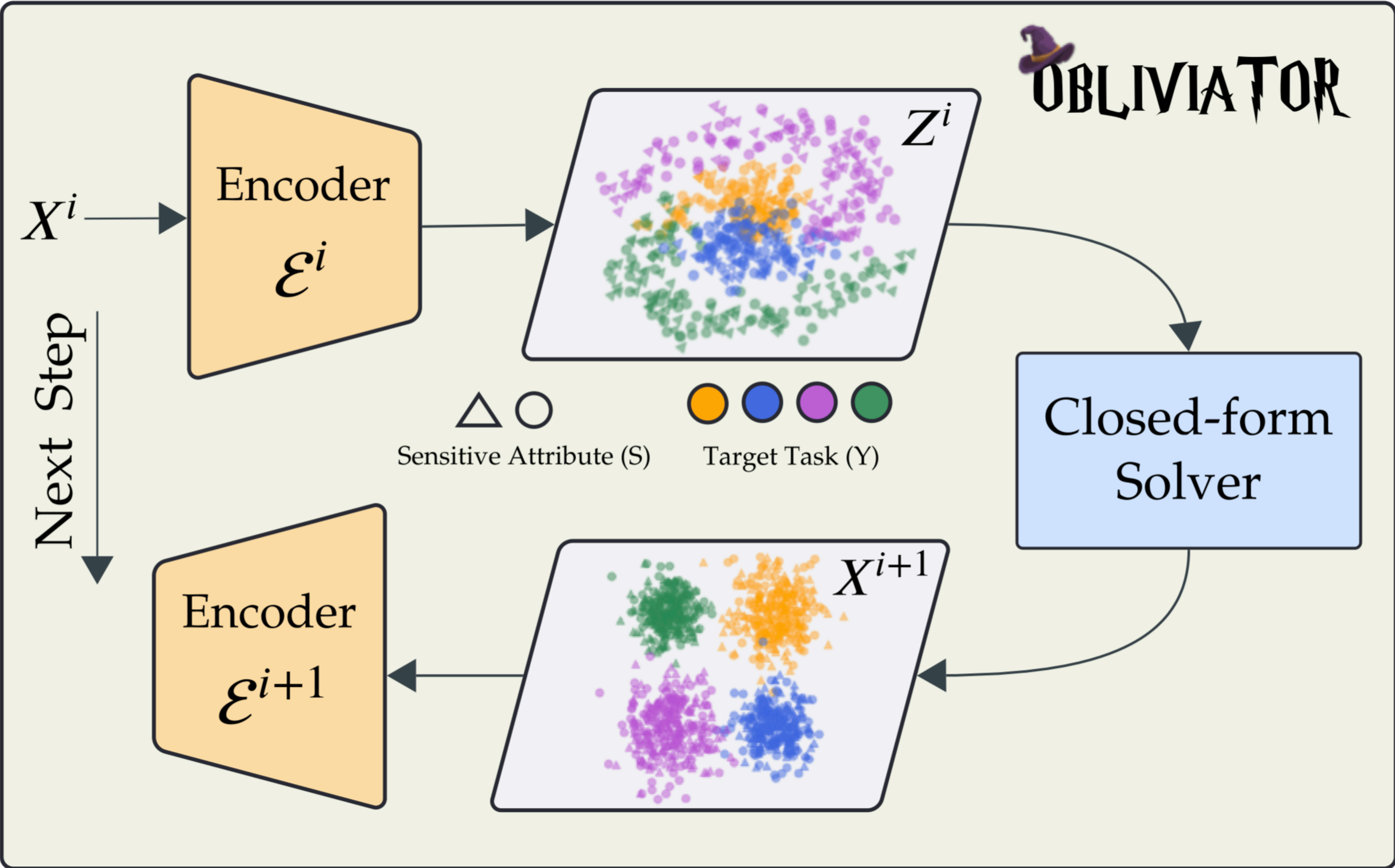
Single-Shot Optimization Leads to **Poor** Solutions

To Make Erasure **Smoother** We Propose an **Iterative** approach

How Do We **Solve** This Optimization ?

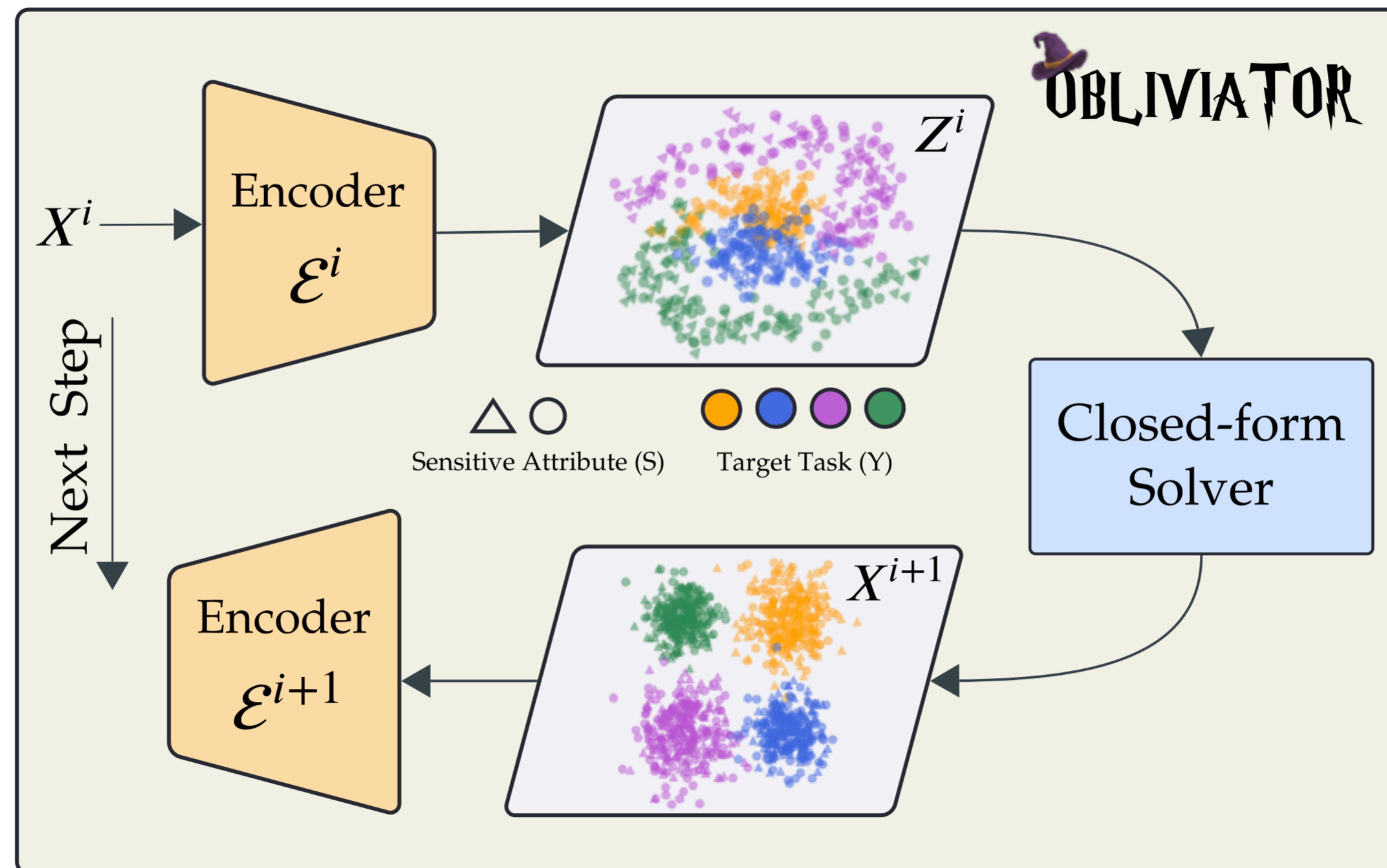


How Do We **Solve** This Optimization ?



Step One

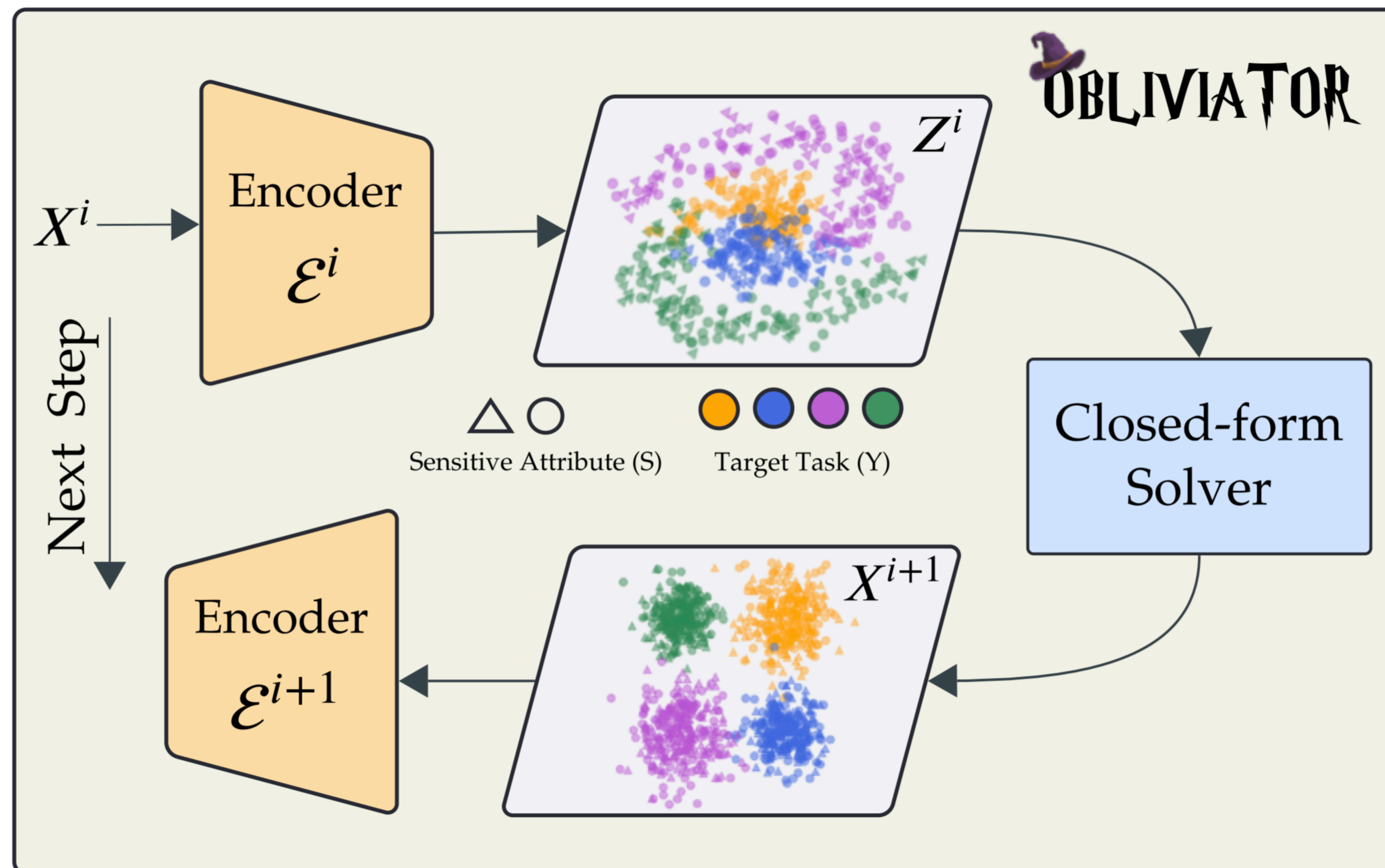
How Do We **Solve** This Optimization ?



Step One

Erase and **Preserve** via guarding function $Z_{\theta}^i = \varepsilon^i(X^i; \theta^i)$

How Do We **Solve** This Optimization ?

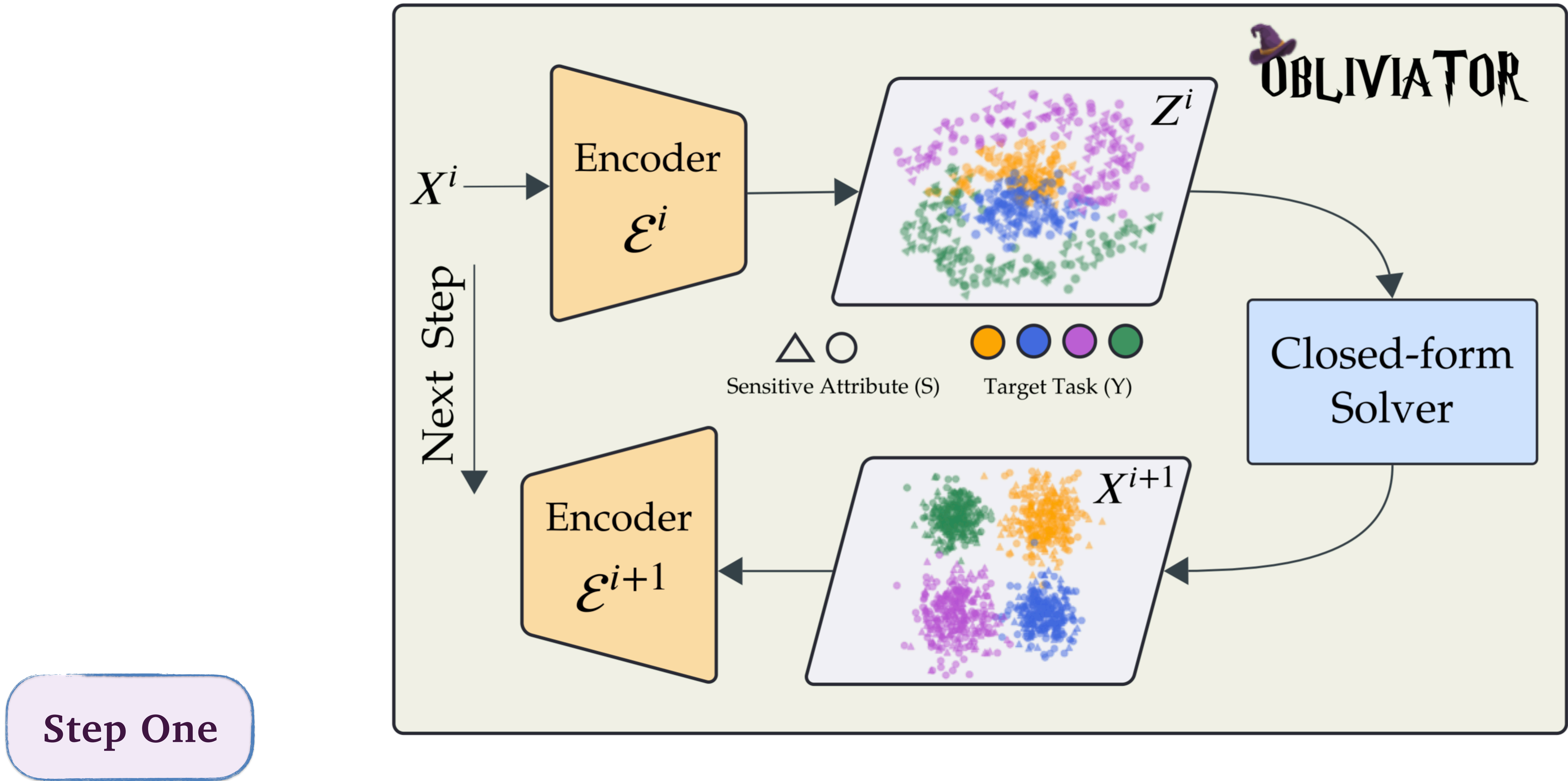


Step One

Erase and **Preserve** via guarding function $Z_\theta^i = \mathcal{E}^i(X^i; \theta^i)$

$$\inf_{\theta} \quad \text{HSIC}(Z_\theta, S) - [\text{HSIC}(Z_\theta^i, Y) + \text{HSIC}(Z_\theta^i, X) + \text{HSIC}(Z_\theta^i, X^i)]$$

How Do We **Solve** This Optimization ?



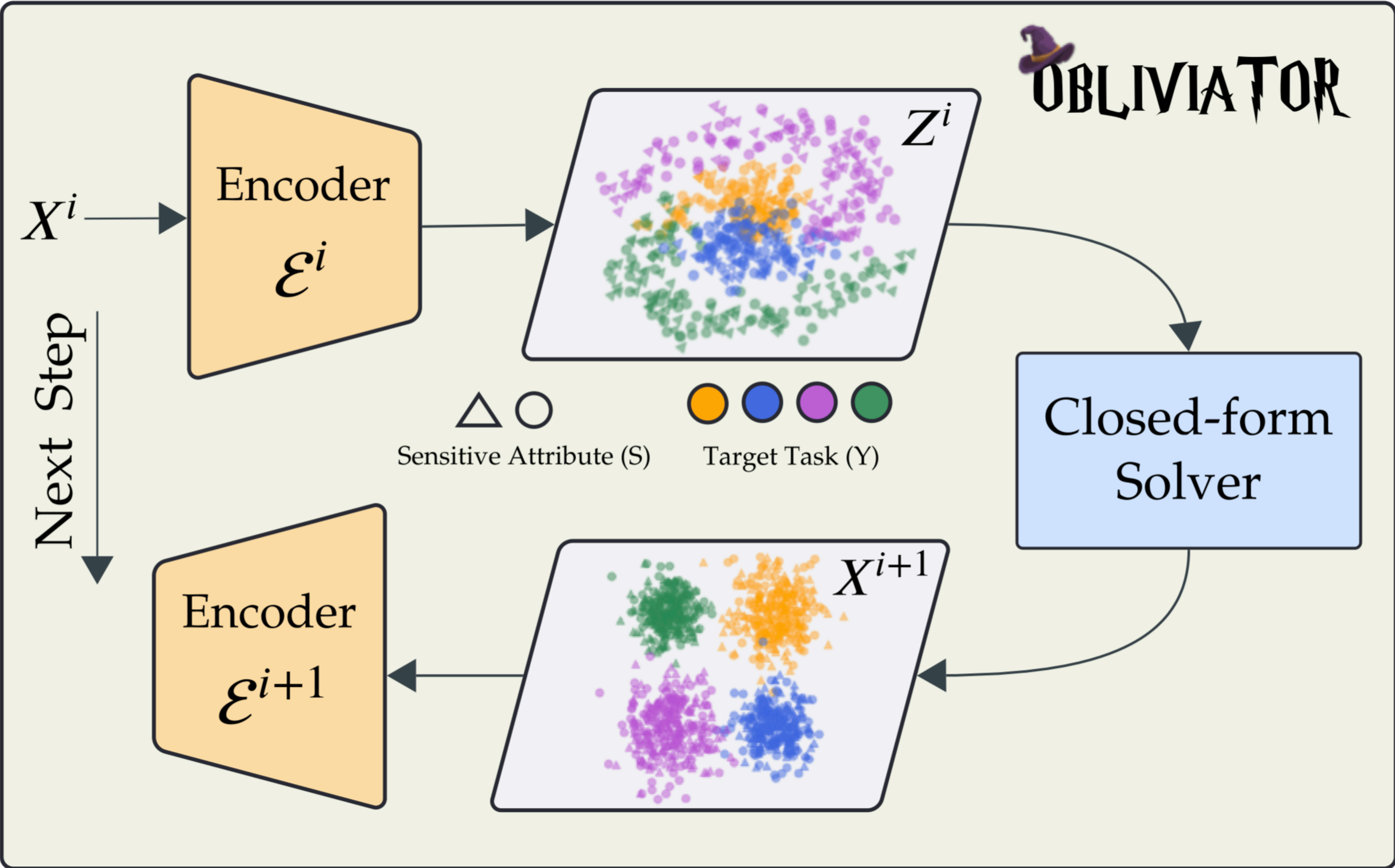
Erase and **Preserve** via guarding function $Z_{\theta}^i = \varepsilon^i(X^i; \theta^i)$

\inf_{θ}

$$\text{HSIC}(Z_{\theta}, S) - [\text{HSIC}(Z_{\theta}^i, Y) + \text{HSIC}(Z_{\theta}^i, X) + \text{HSIC}(Z_{\theta}^i, X^i)]$$

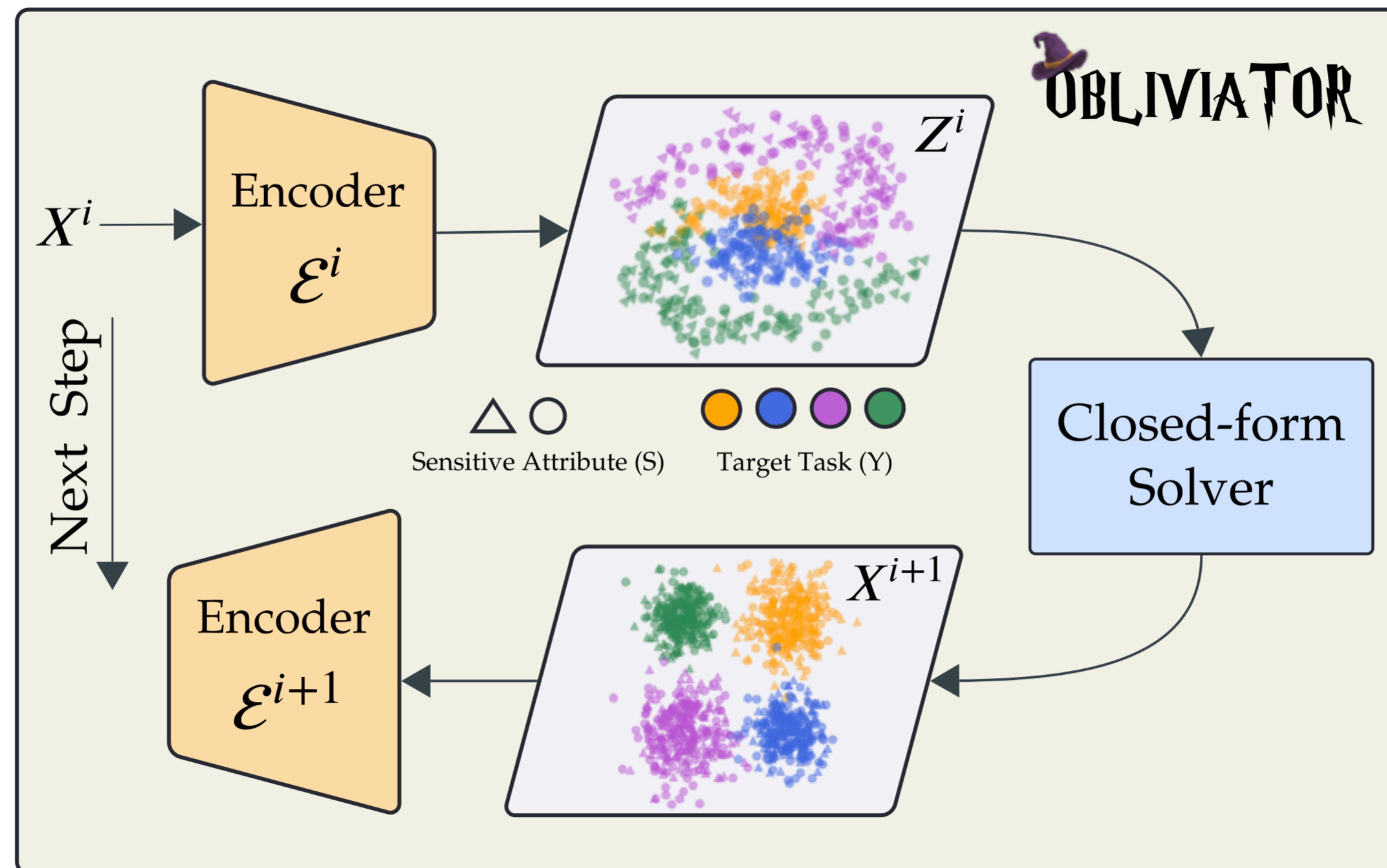
More Information-Preserving Loss

How Do We **Solve** This Optimization ?



Step Two

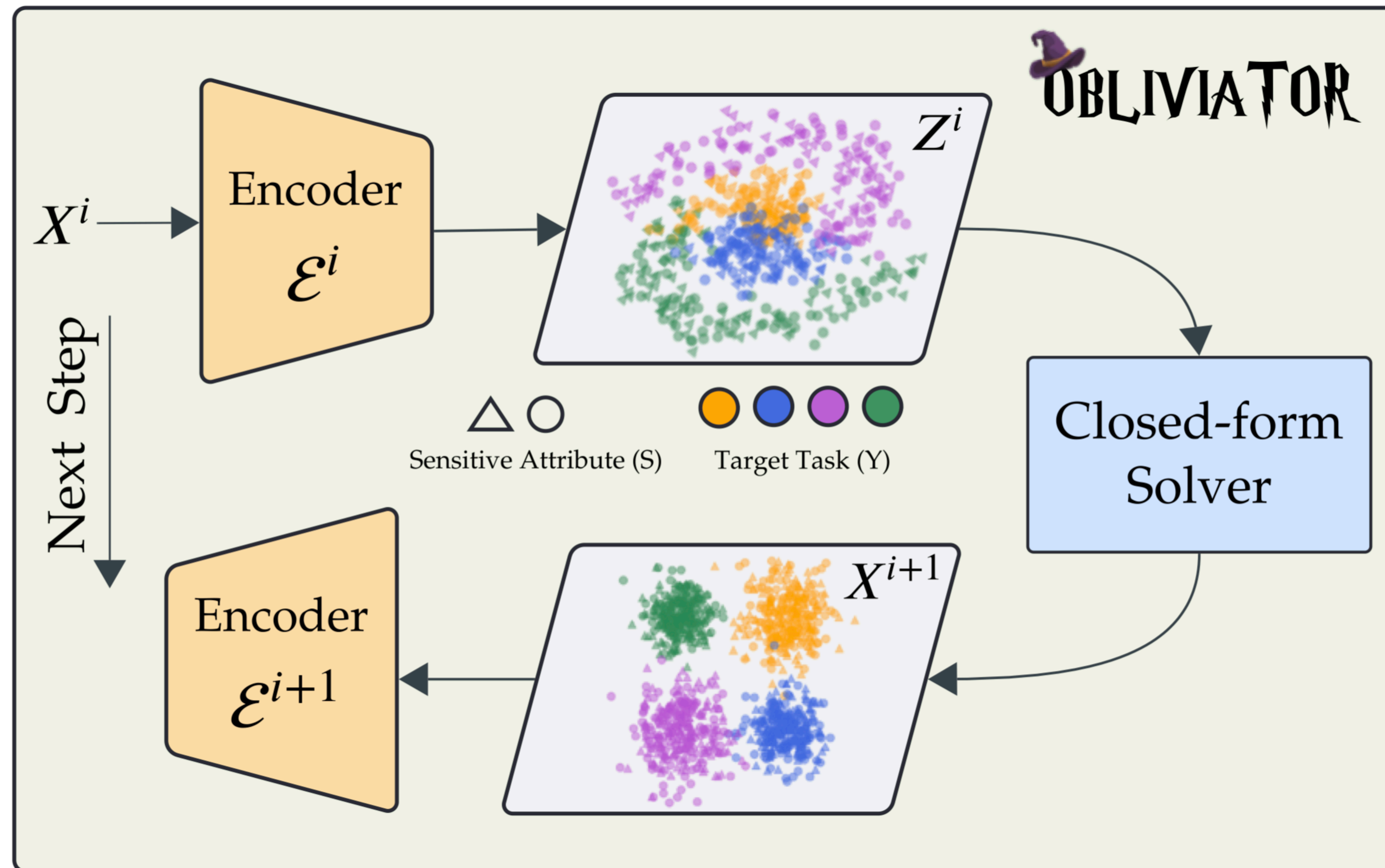
How Do We **Solve** This Optimization ?



Step Two

Separate **Useful Information** from **Unwanted Concepts**:

How Do We **Solve** This Optimization ?



Step Two

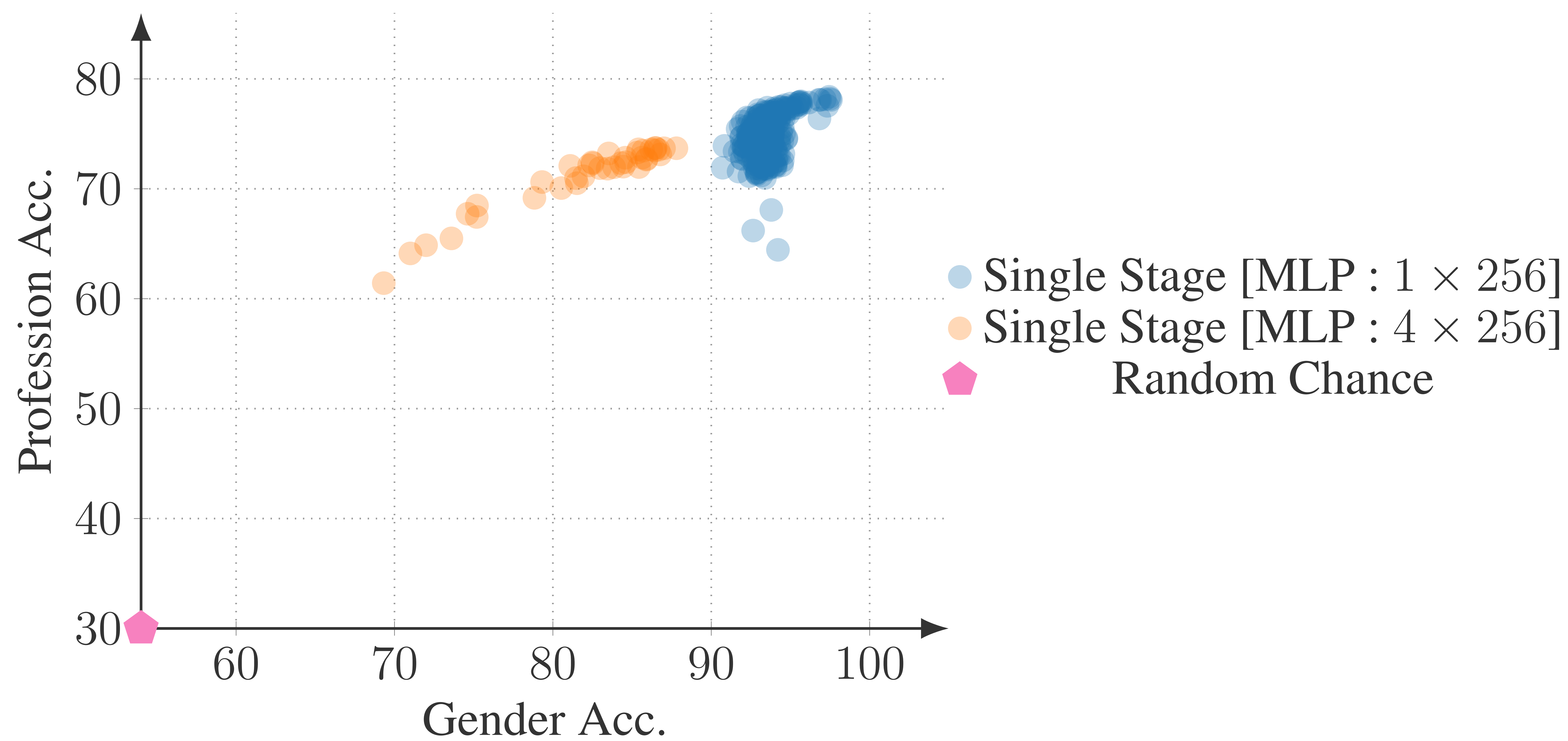
Separate **Useful Information** from **Unwanted Concepts**:

$$\sup_{\{g_a \in \mathcal{G}_a\}} \sup_{f \in \mathcal{F}} \mathbf{Cov}^2(f(Z_\theta^i), g_y(Y)) + \mathbf{Cov}^2(f(Z_\theta^i), g_{x^i}(X^i)) + \mathbf{Cov}^2(f(Z_\theta^i), g_x(X))$$

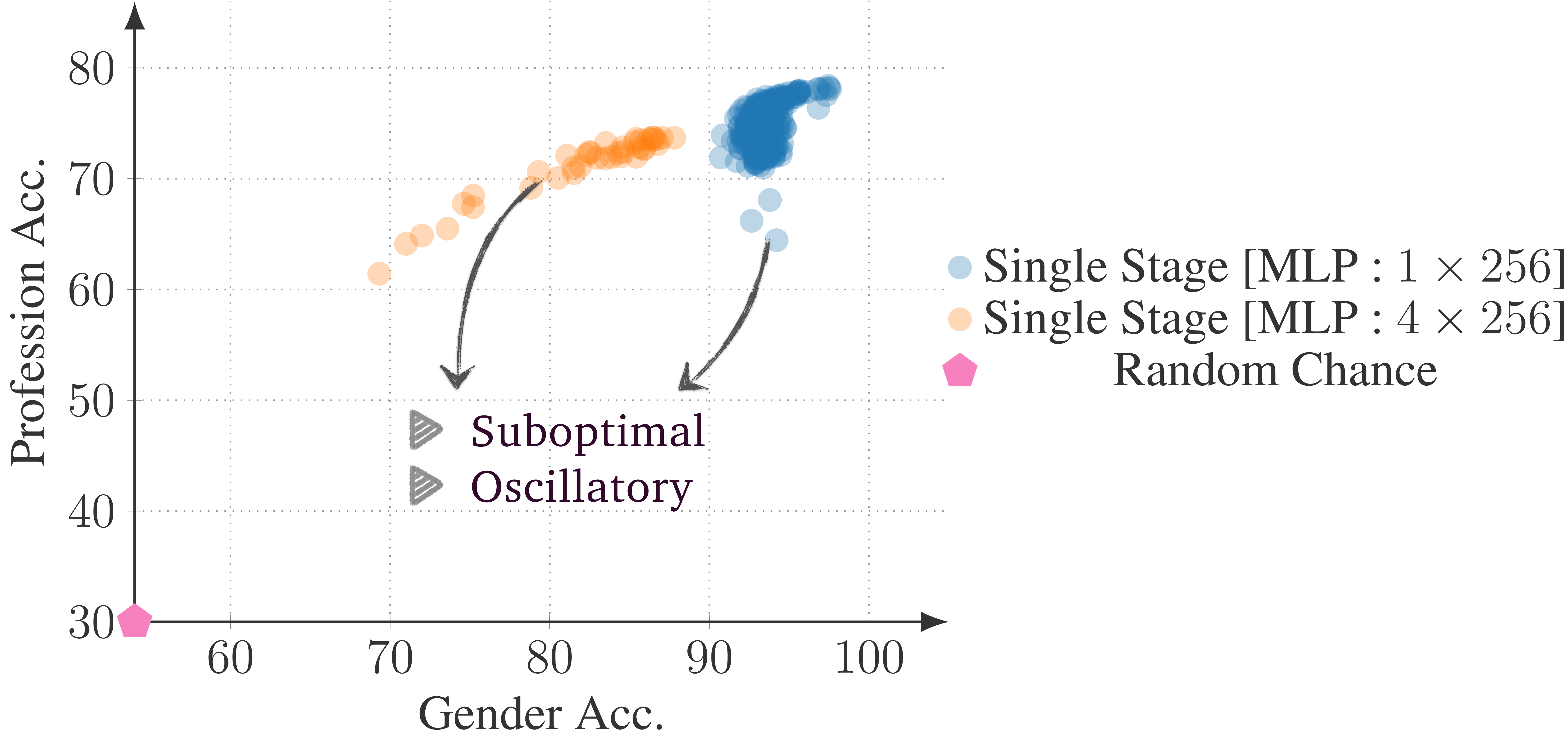
$$s.t. \quad \sup_{g_s \in \mathcal{G}_s} \mathbf{Cov}(f(Z_\theta^i), g_s(S)) = 0 \quad \|f\|_{\mathcal{F}} = \|g_a\|_{\mathcal{G}_a} = 1$$

Single Stage Vs. **Iterative** Erasure

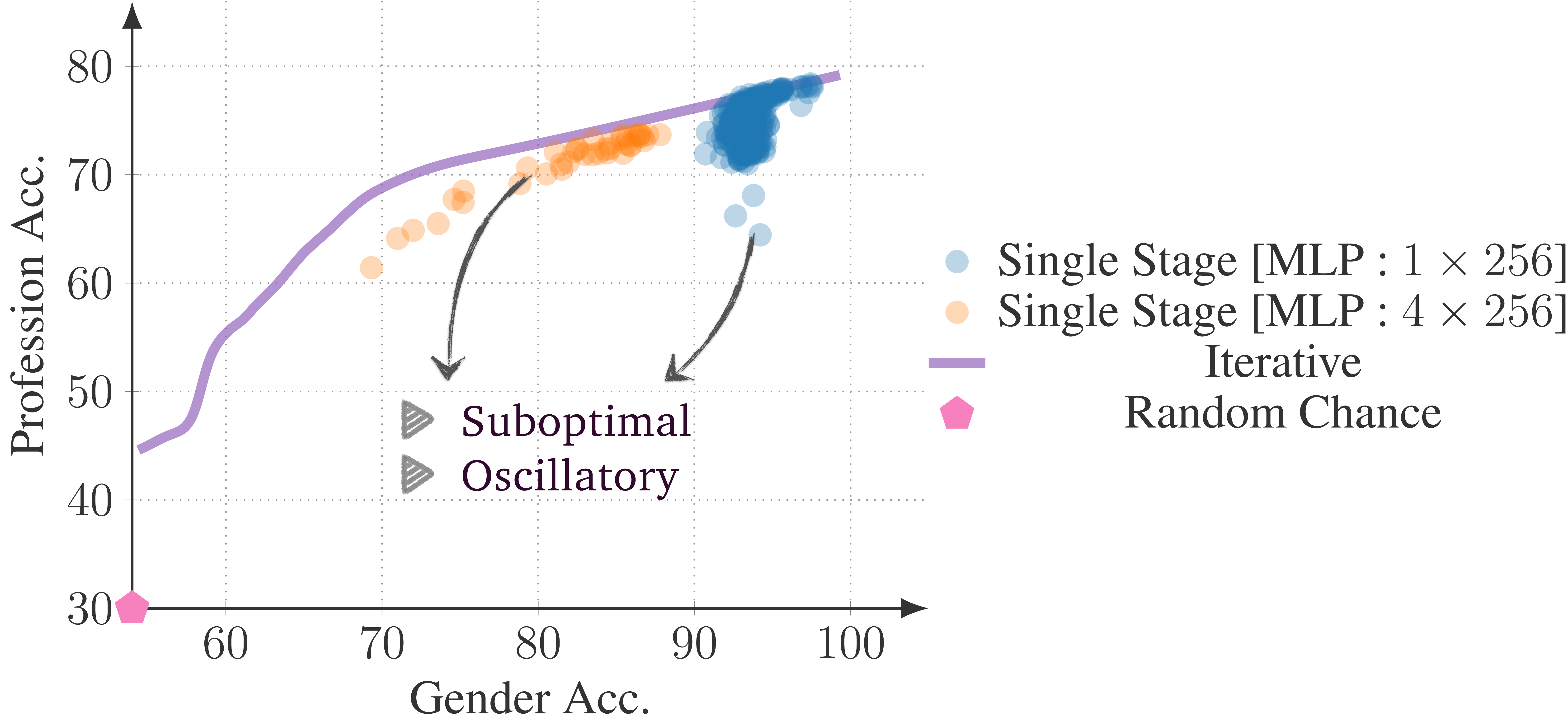
Single Stage Vs. Iterative Erasure



Single Stage Vs. **Iterative** Erasure

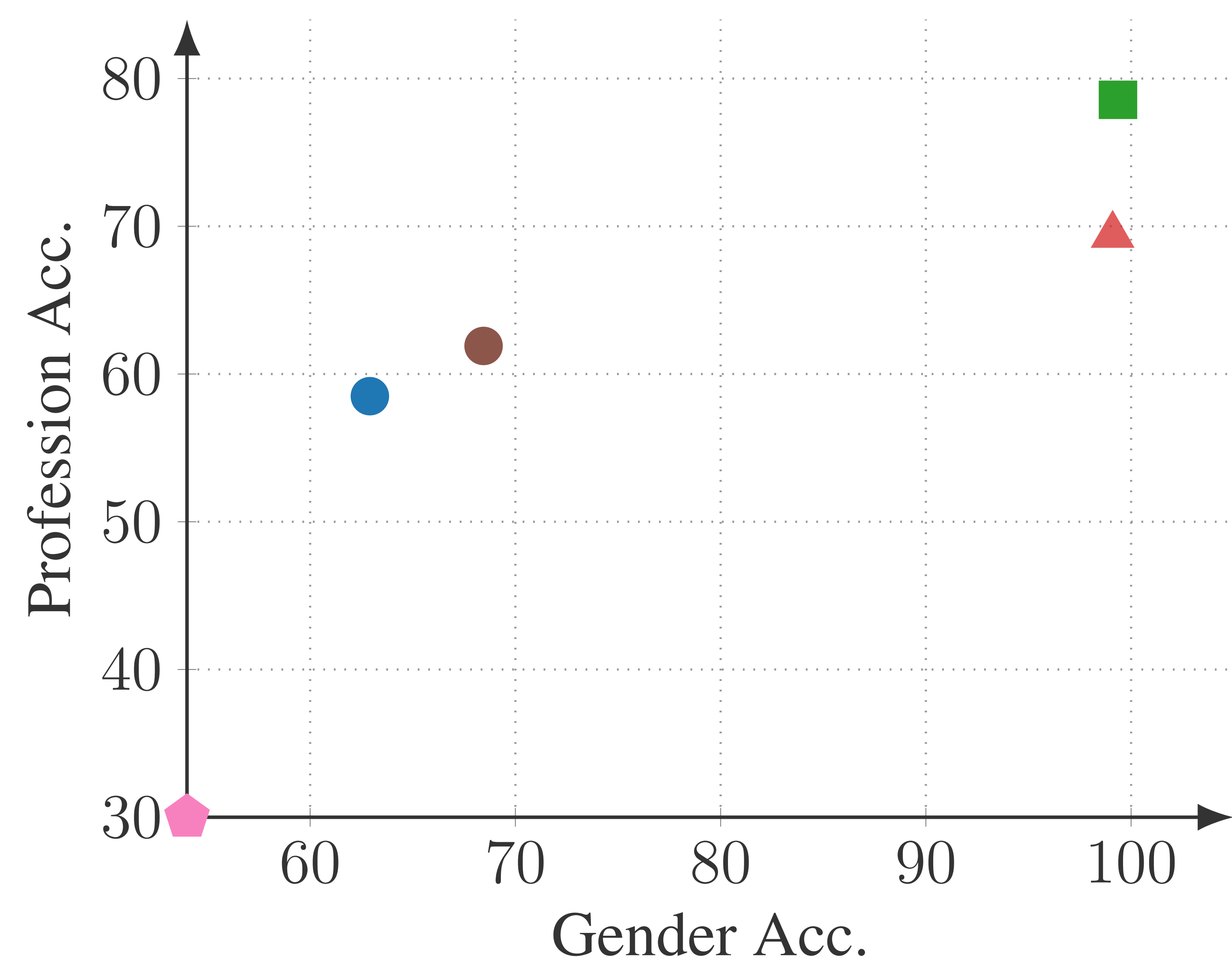


Single Stage Vs. **Iterative** Erasure



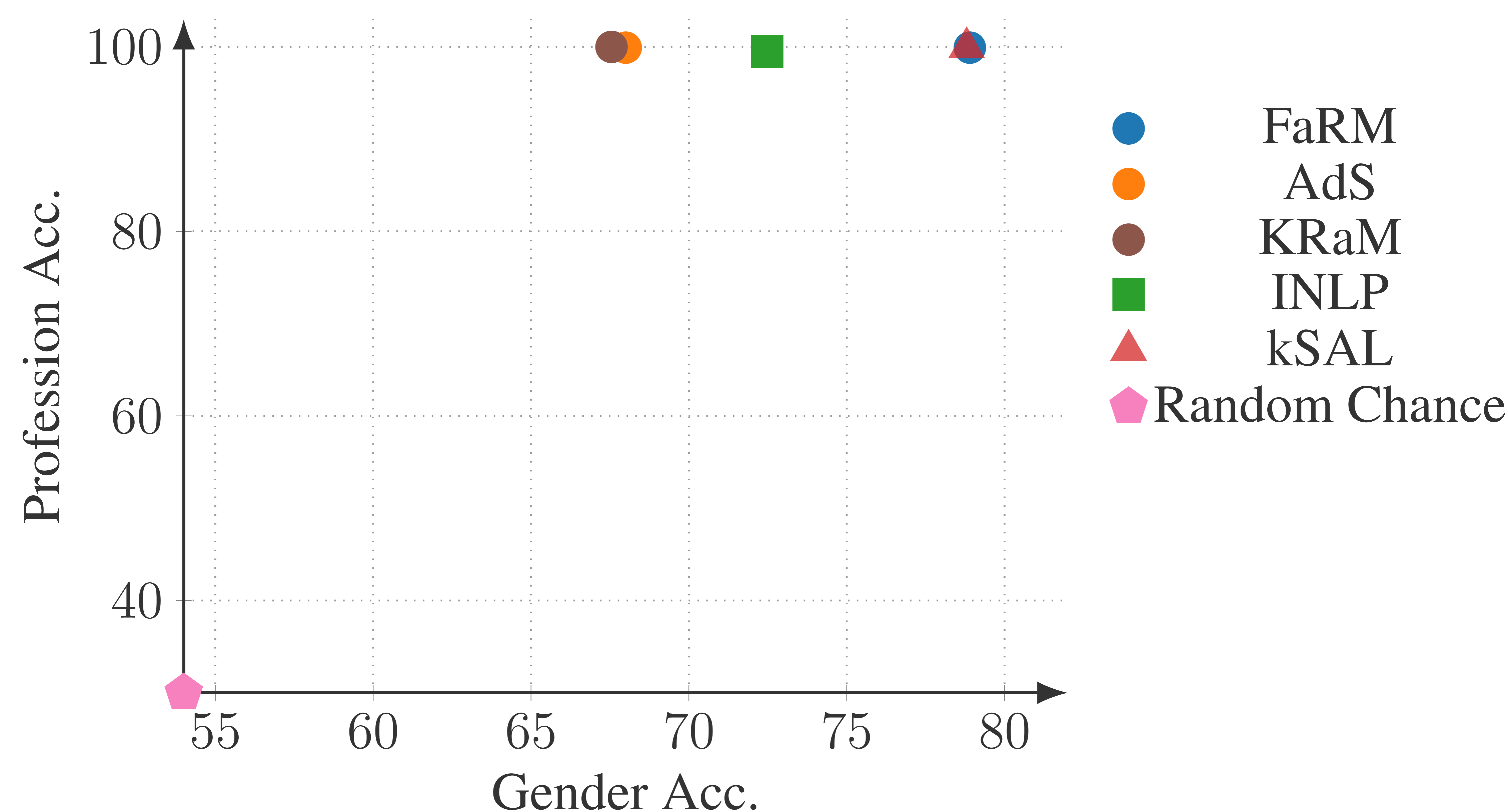
What does Obliviator **Reveal about **Concept Erasure** ?**

Obliviator Reveals : Achieving **Nonlinear** Guardedness



BIAS IN BIOS

Representation: **Frozen**



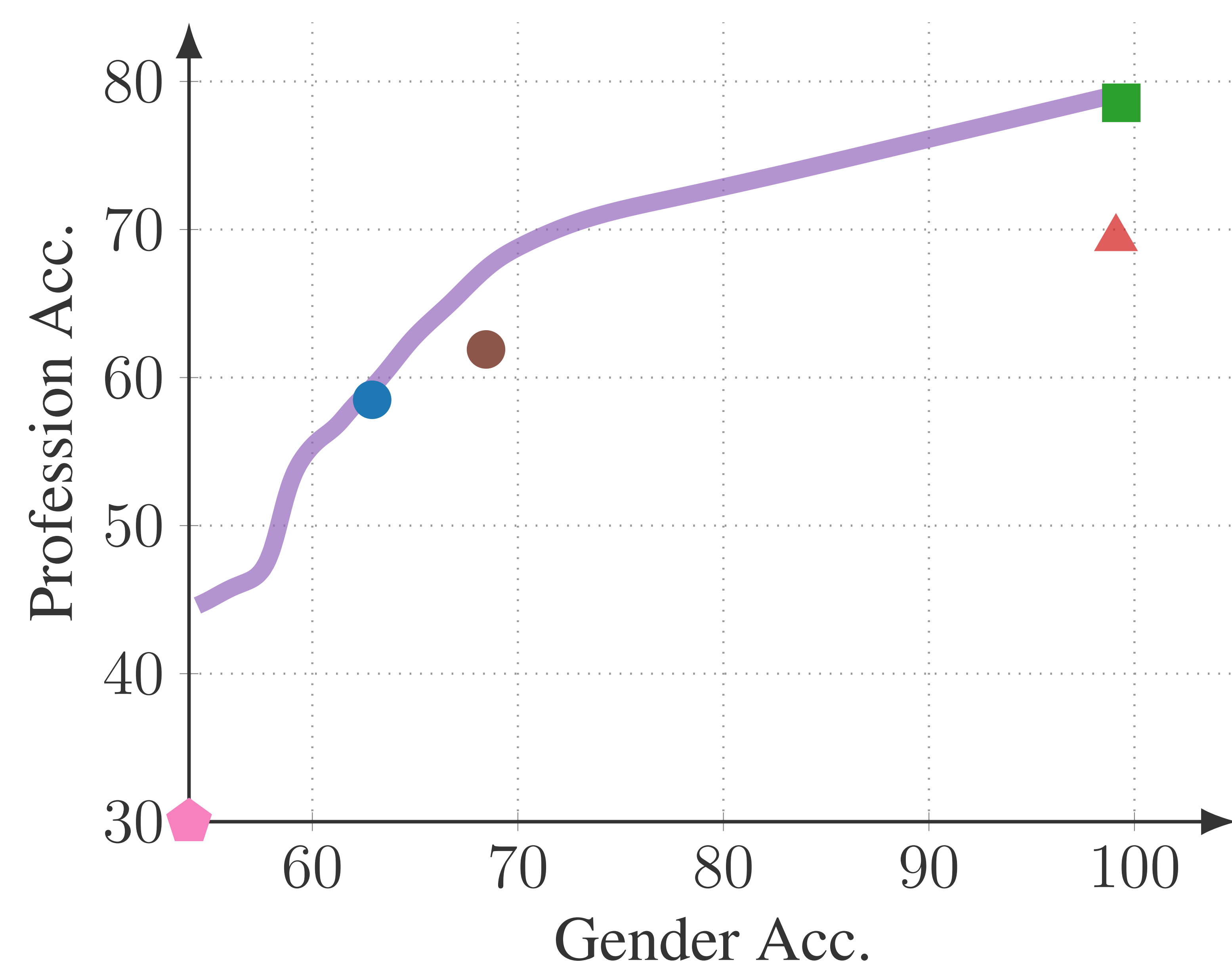
BIAS IN BIOS

Representation: **Finetuned**

PLM : **BERT**

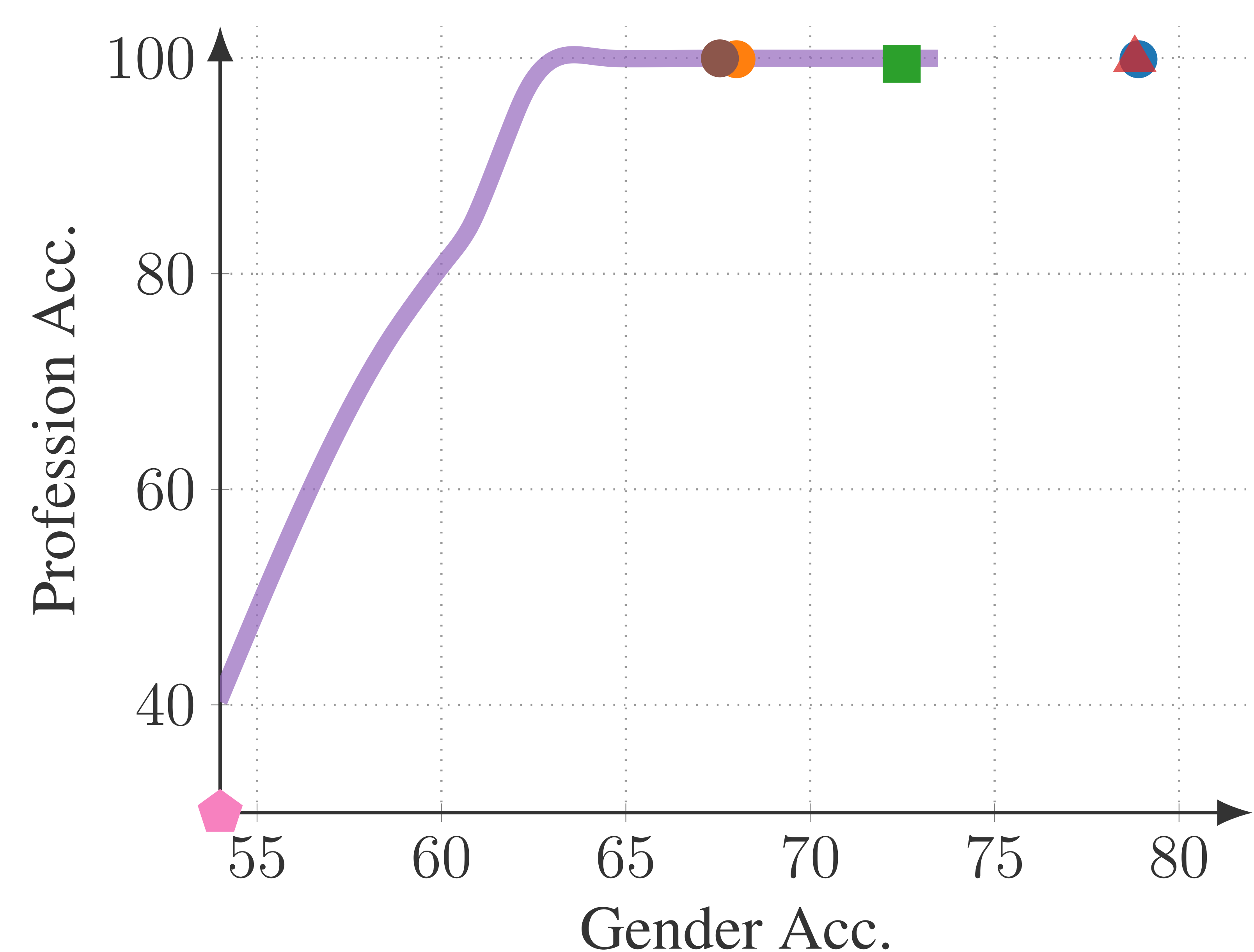
Utility : Profession
Unwanted : Gender

Obliviator Reveals : Achieving **Nonlinear** Guardedness



BIAS IN BIOS

Representation: **Frozen**



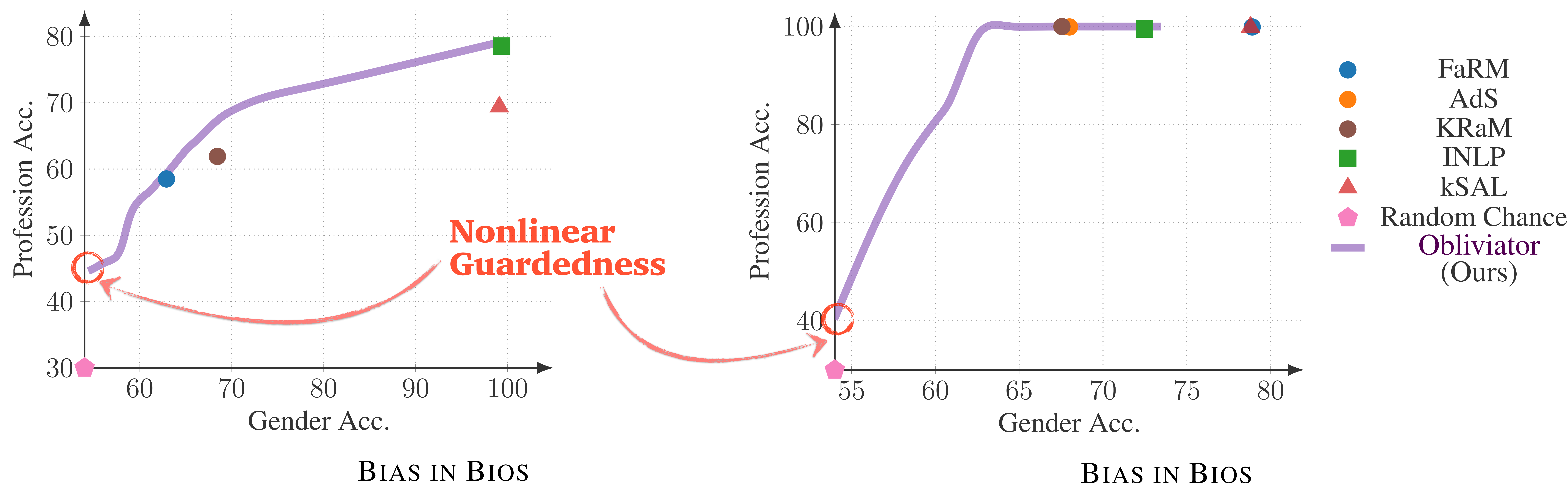
BIAS IN BIOS

Representation: **Finetuned**

PLM : **BERT**

Utility : Profession
Unwanted : Gender

Obliviator Reveals : Achieving **Nonlinear** Guardedness



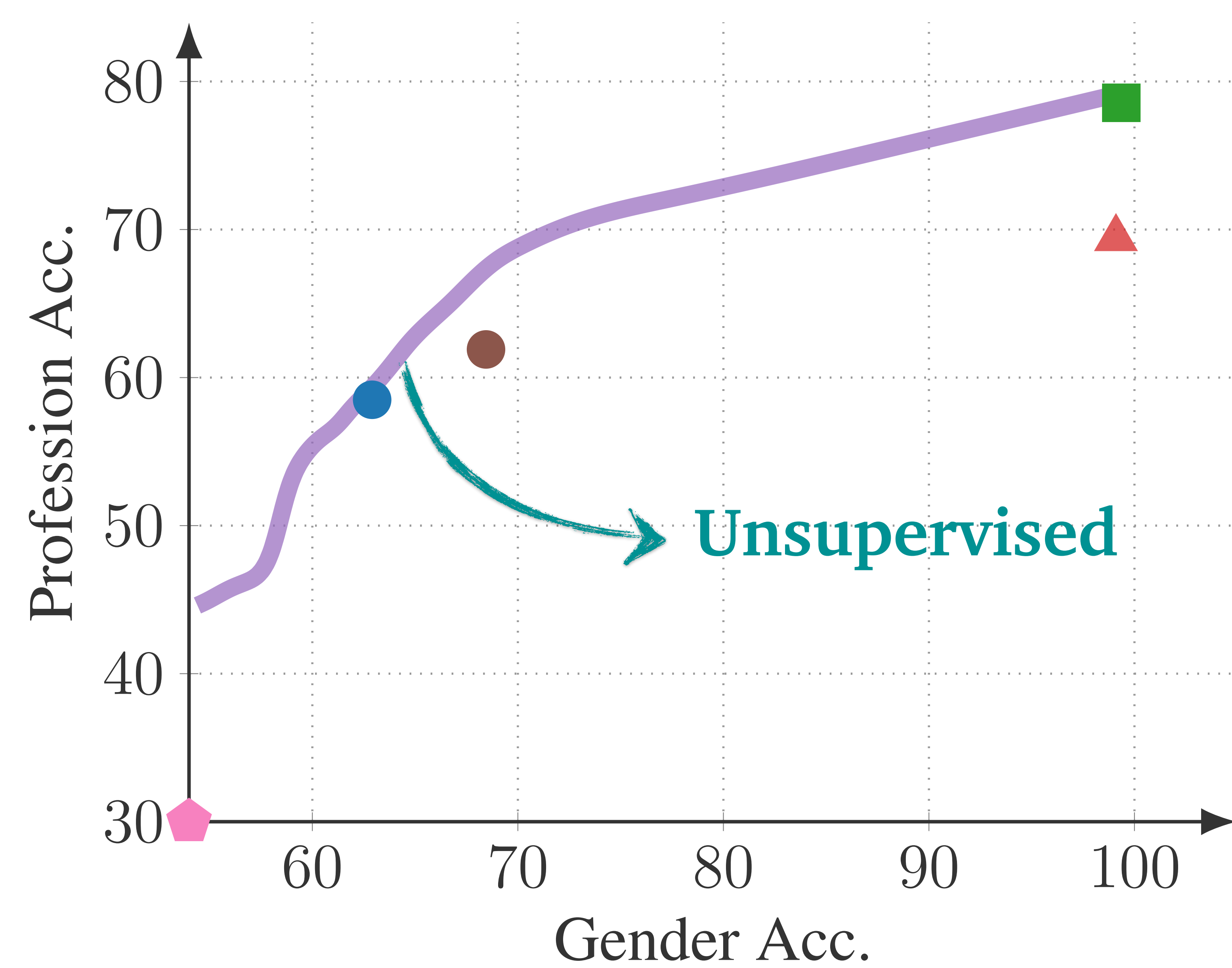
Representation: **Frozen**

Representation: **Finetuned**

PLM : **BERT**

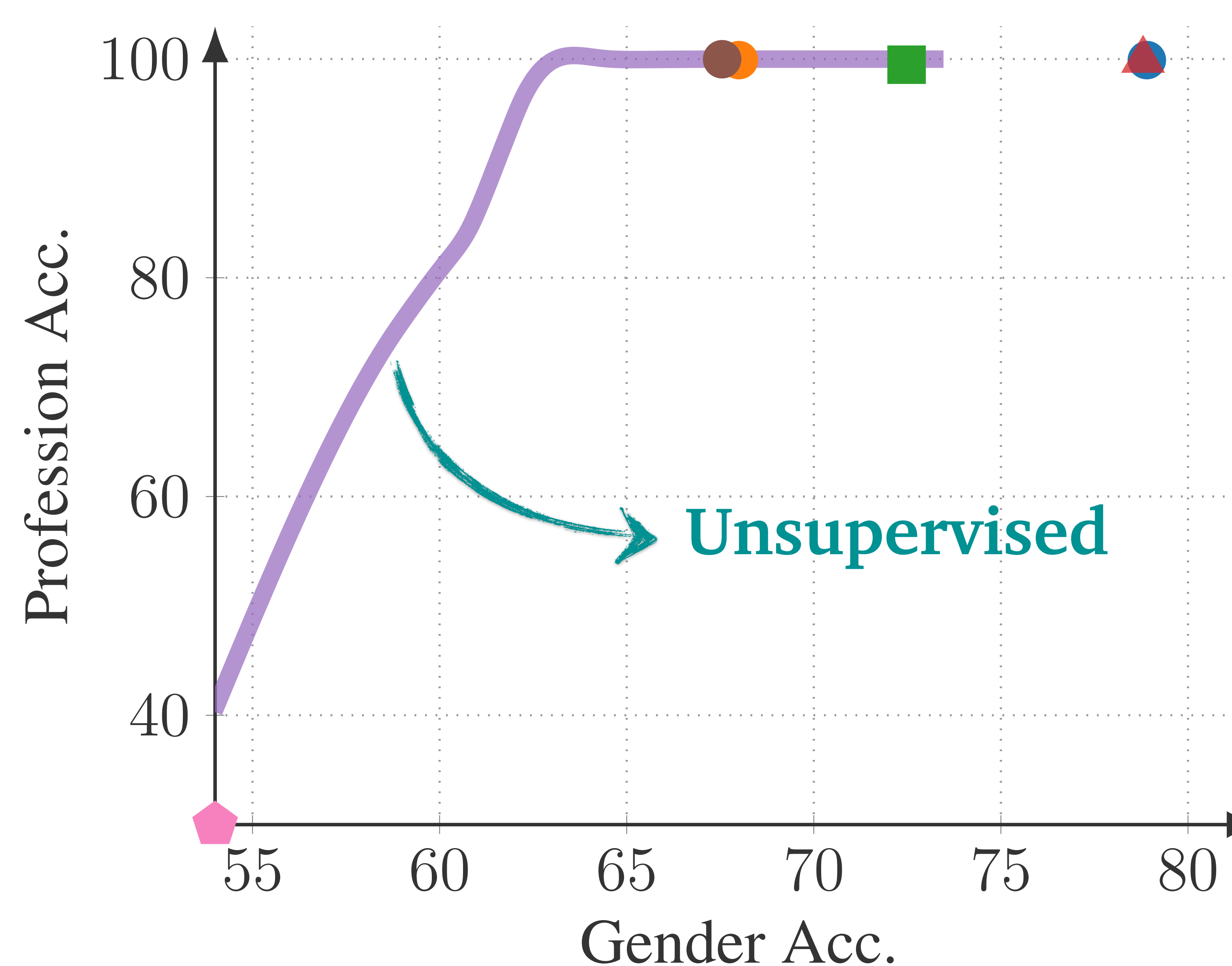
Utility : Profession
Unwanted : Gender

Obliviator Reveals : Effect of **Target Task Labels** on Erasure



BIAS IN BIOS

Representation: **Frozen**



BIAS IN BIOS

Representation: **Finetuned**

PLM : **BERT**

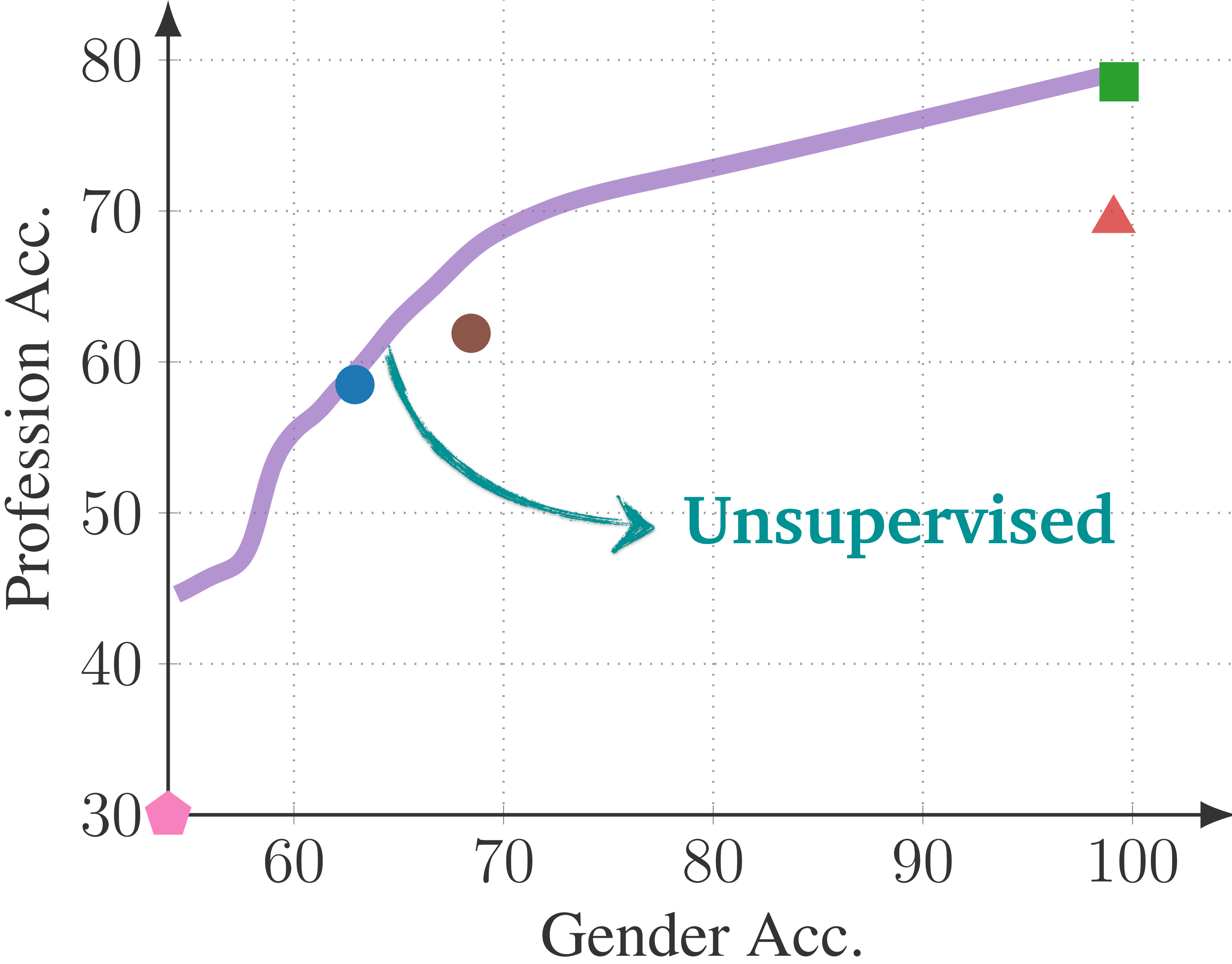
Utility : Profession
Unwanted : Gender

- FaRM
- AdS
- KRaM
- INLP
- kSAL
- Random Chance
- Obliviator (Ours)

Obliviator Reveals : Effect of **Target Task Labels** on Erasure

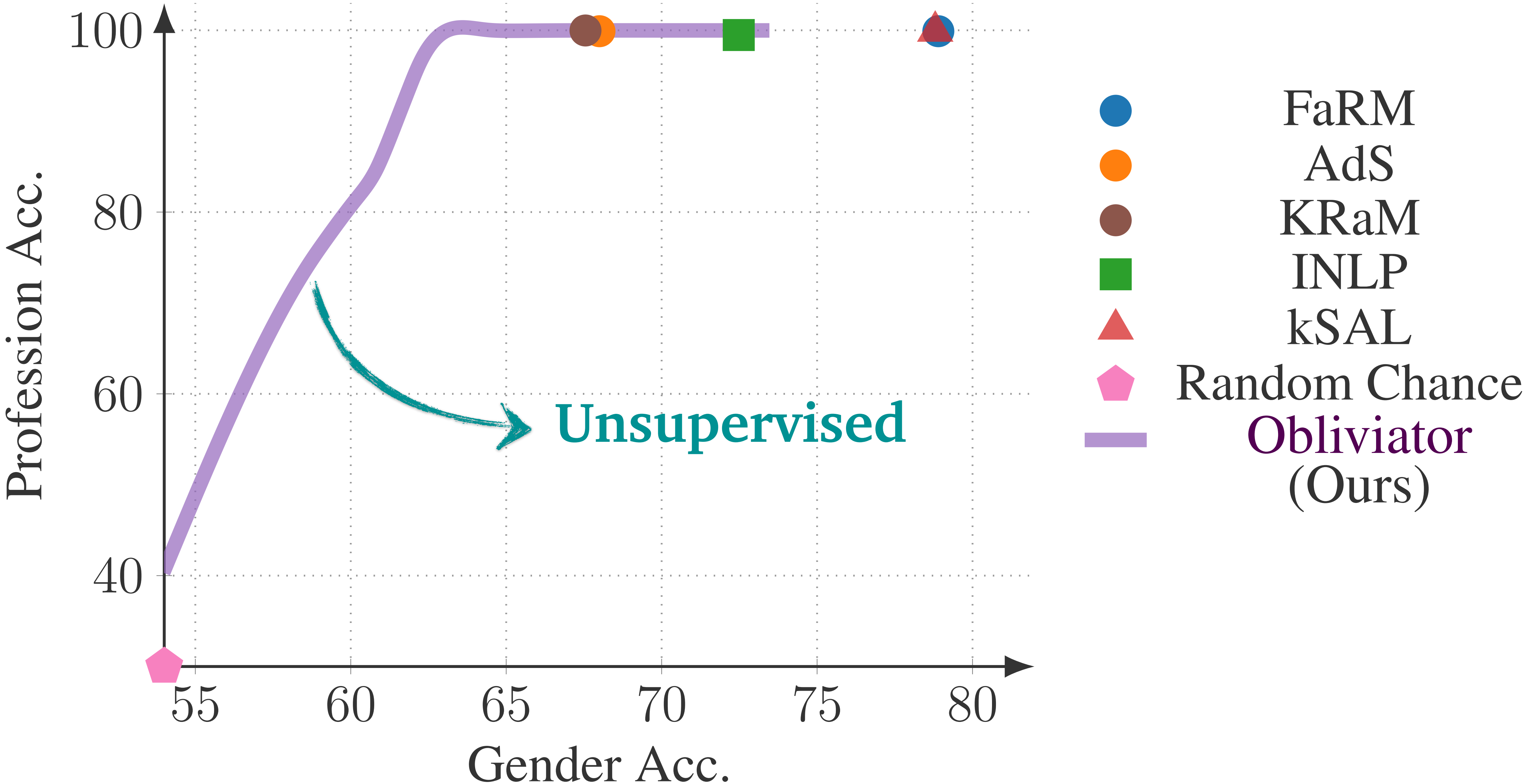
Unsupervised

$$\inf_{\theta} \text{HSIC}(Z_{\theta}, S) - [\text{HSIC}(Z_{\theta}^i, Y) + \text{HSIC}(Z_{\theta}^i, X) + \text{HSIC}(Z_{\theta}^i, X^i)]$$



BIAS IN BIOS

Representation: **Frozen**



BIAS IN BIOS

Representation: **Finetuned**

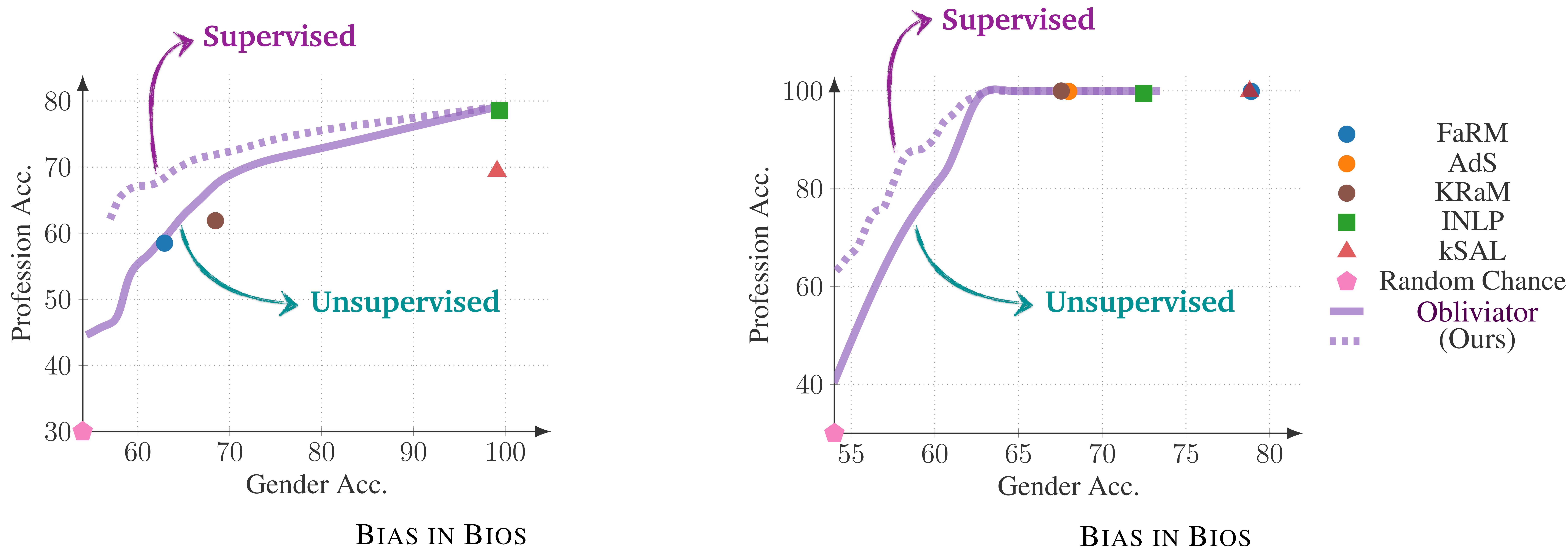
PLM : **BERT**

Utility : **Profession**
Unwanted : **Gender**

Obliviator Reveals : Effect of **Target Task Labels** on Erasure

Supervised

$$\inf_{\theta} \text{HSIC}(Z_{\theta}, S) - [\text{HSIC}(Z_{\theta}^i, Y) + \text{HSIC}(Z_{\theta}^i, X) + \text{HSIC}(Z_{\theta}^i, X^i)]$$



Representation: **Frozen**

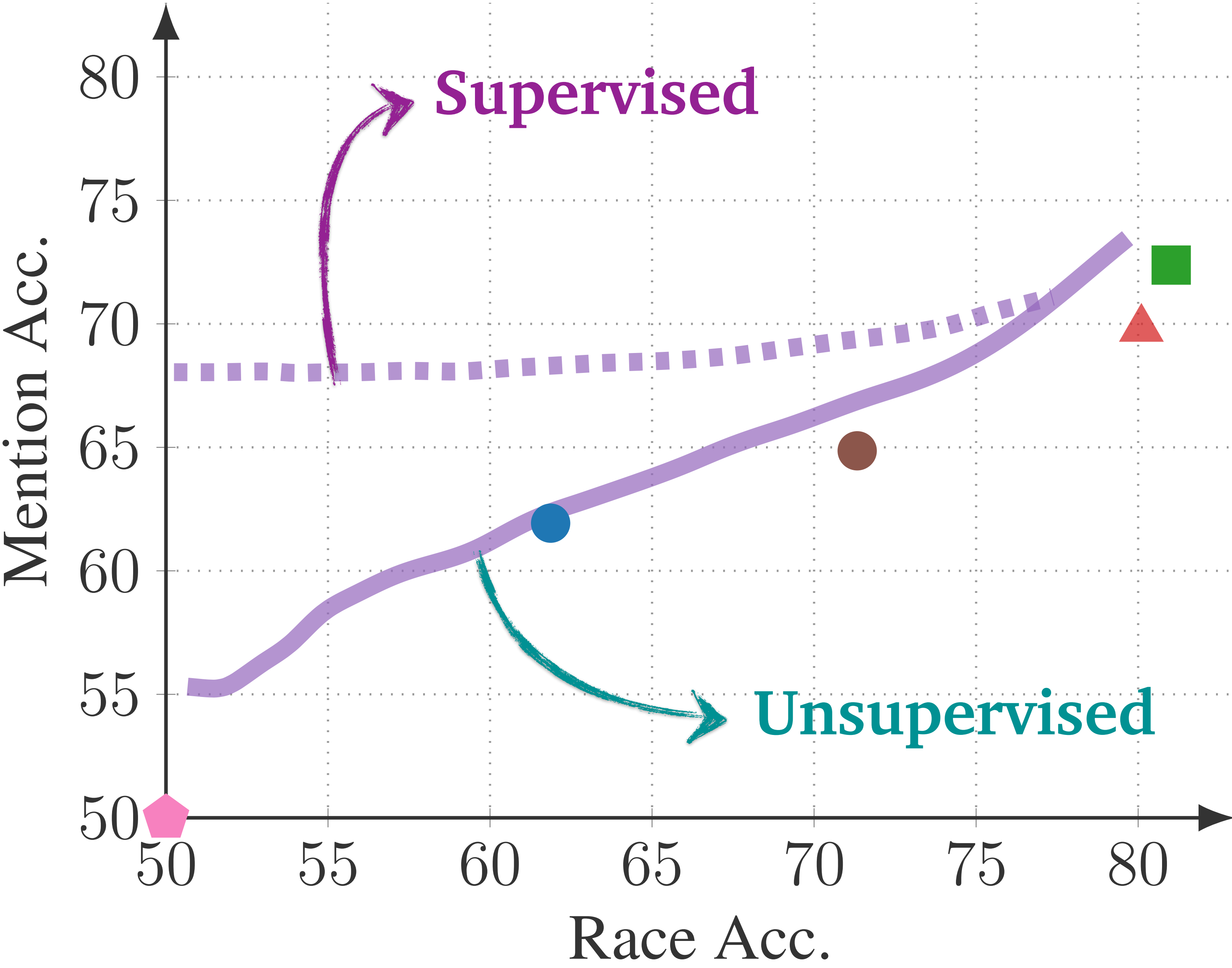
Representation: **Finetuned**

PLM : **BERT**

Utility : Profession
Unwanted : Gender

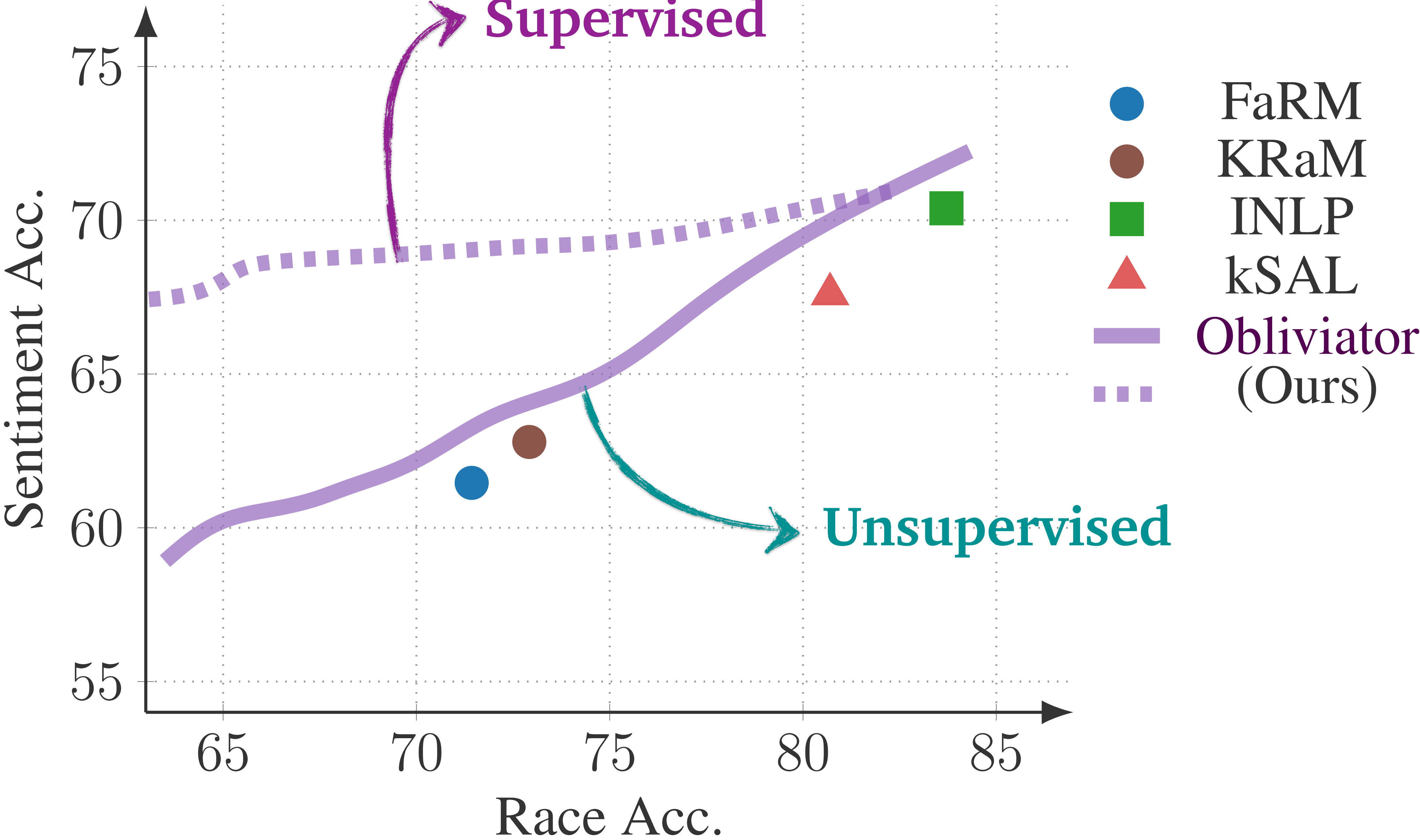
Obliviator Reveals : Effect of **Target Task Labels** on Erasure

Utility : Mention
Unwanted : Race



(a) DIAL-MENTION

Utility : Sentiment
Unwanted : Race



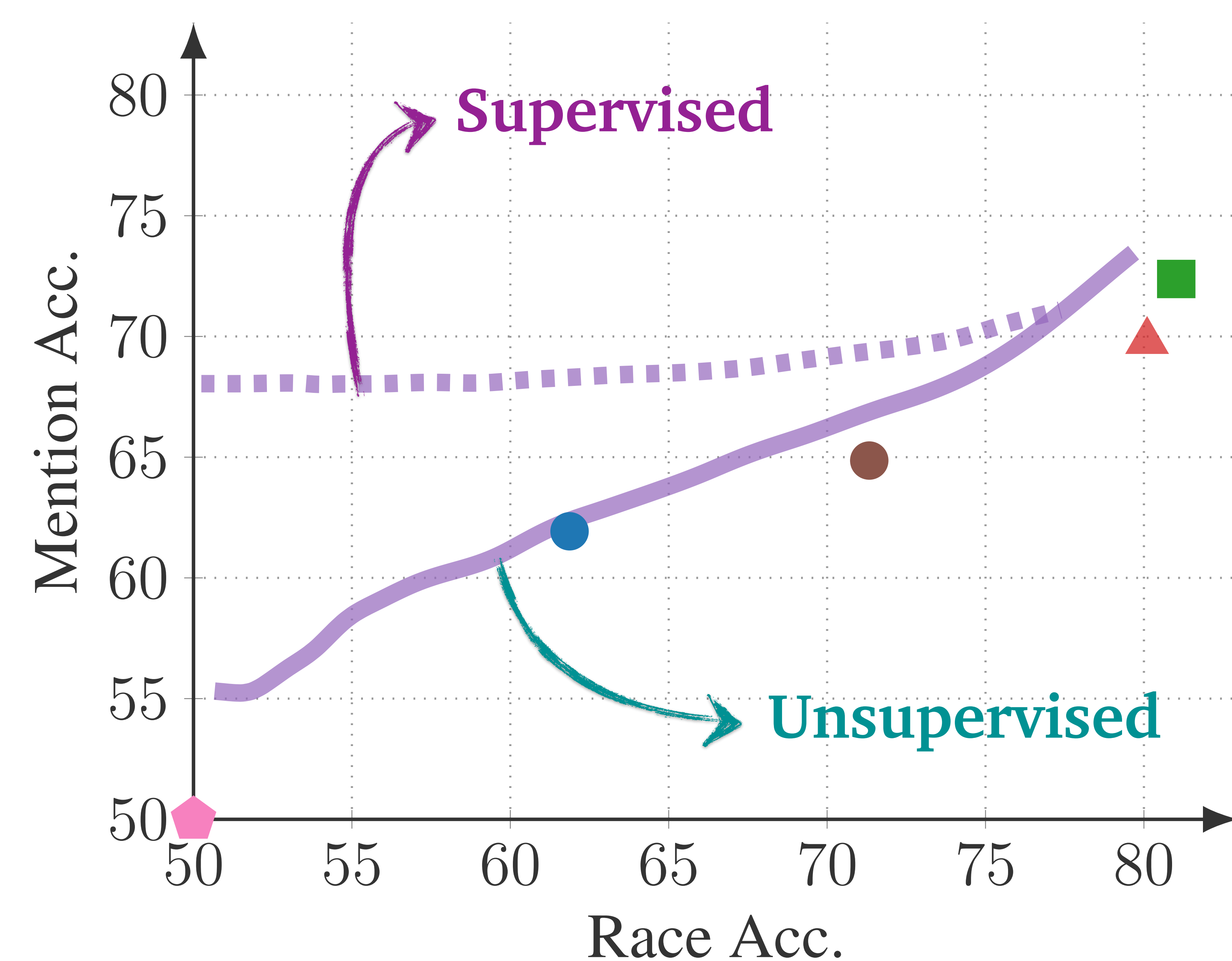
(b) DIAL-SENTIMENT

Representation: **Frozen**

PLM : **BERT**

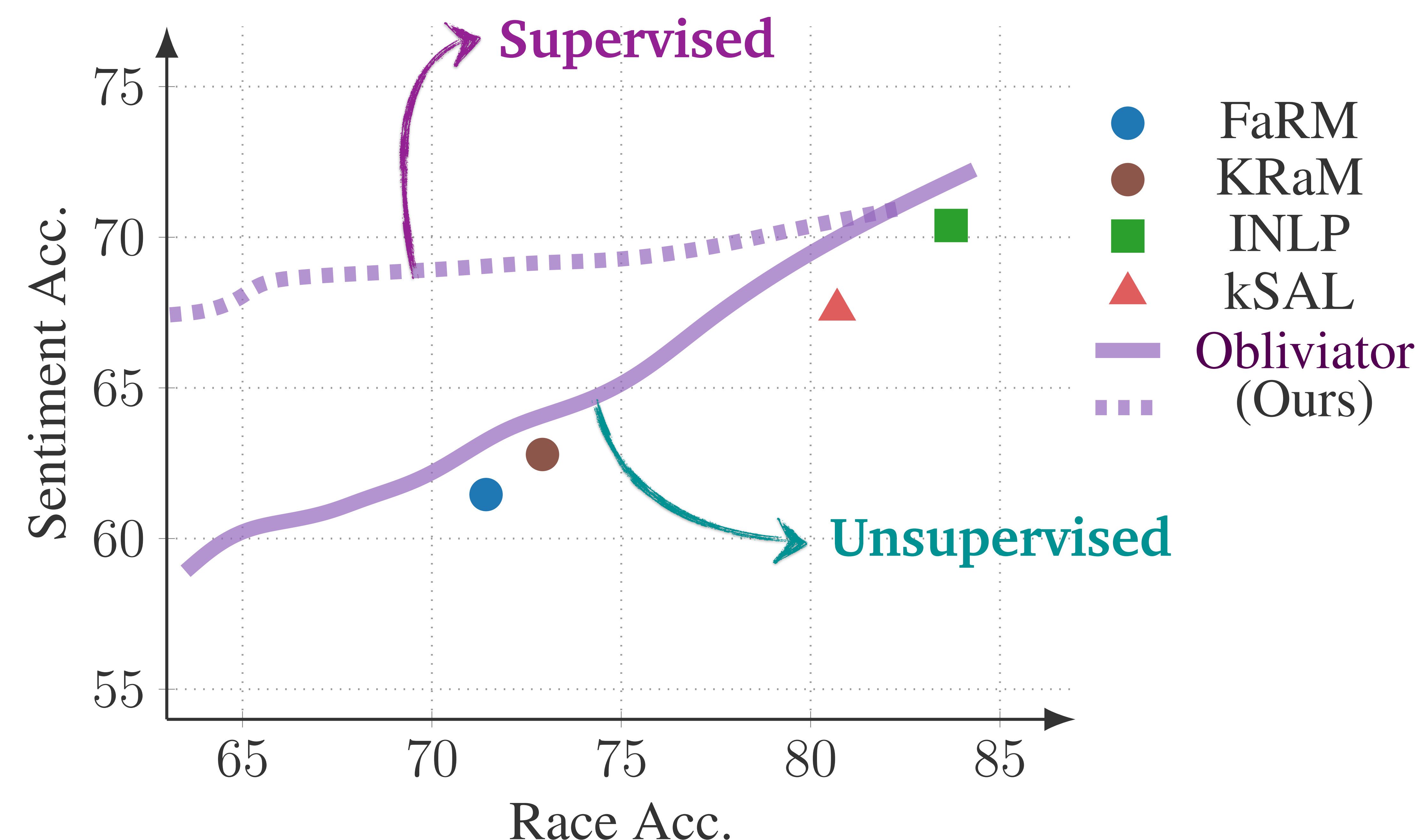
Obliviator Reveals : Effect of **Target Task Labels** on Erasure

Utility : Mention
Unwanted : Race



(a) DIAL-MENTION

Utility : Sentiment
Unwanted : Race

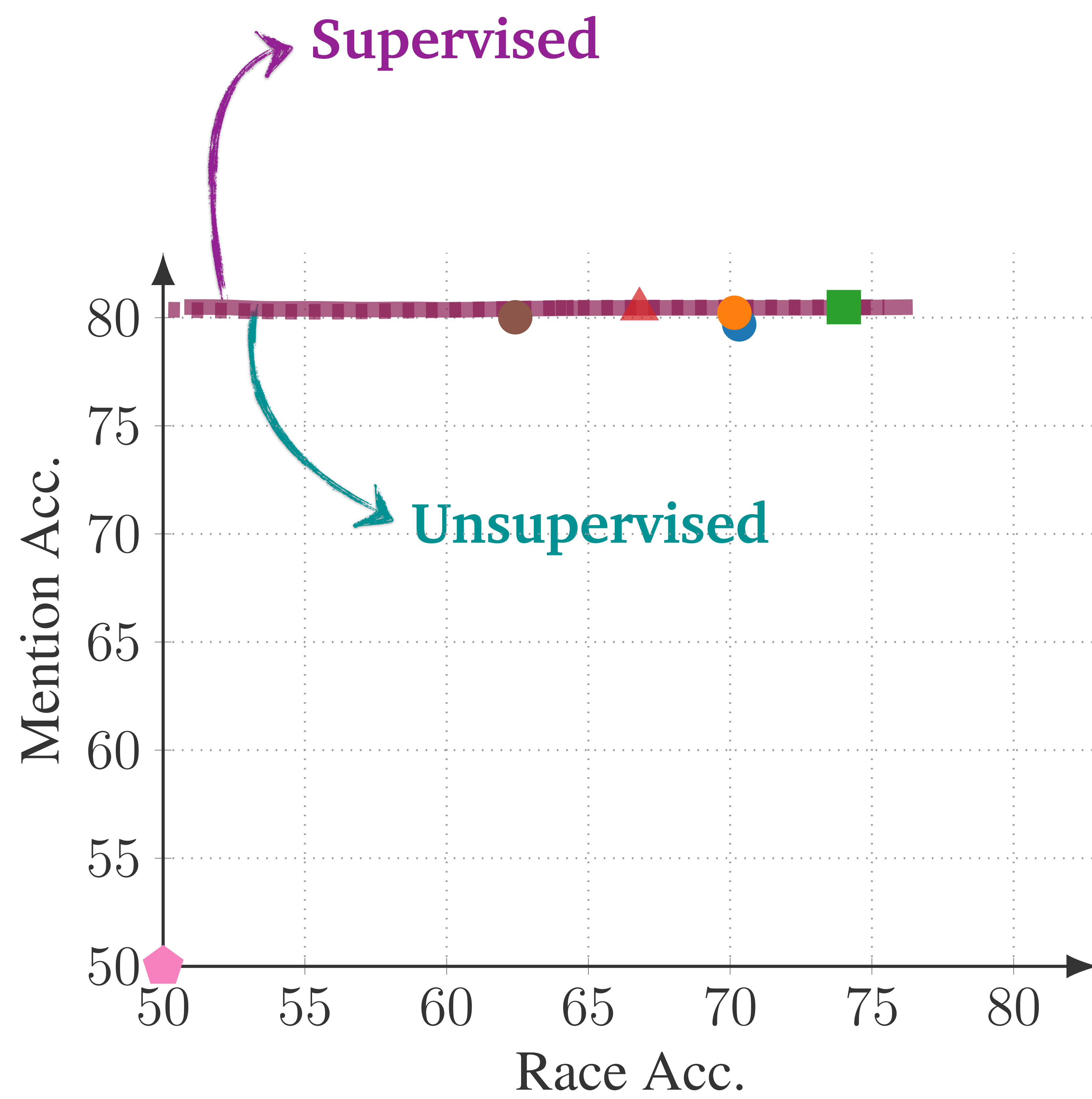


(b) DIAL-SENTIMENT

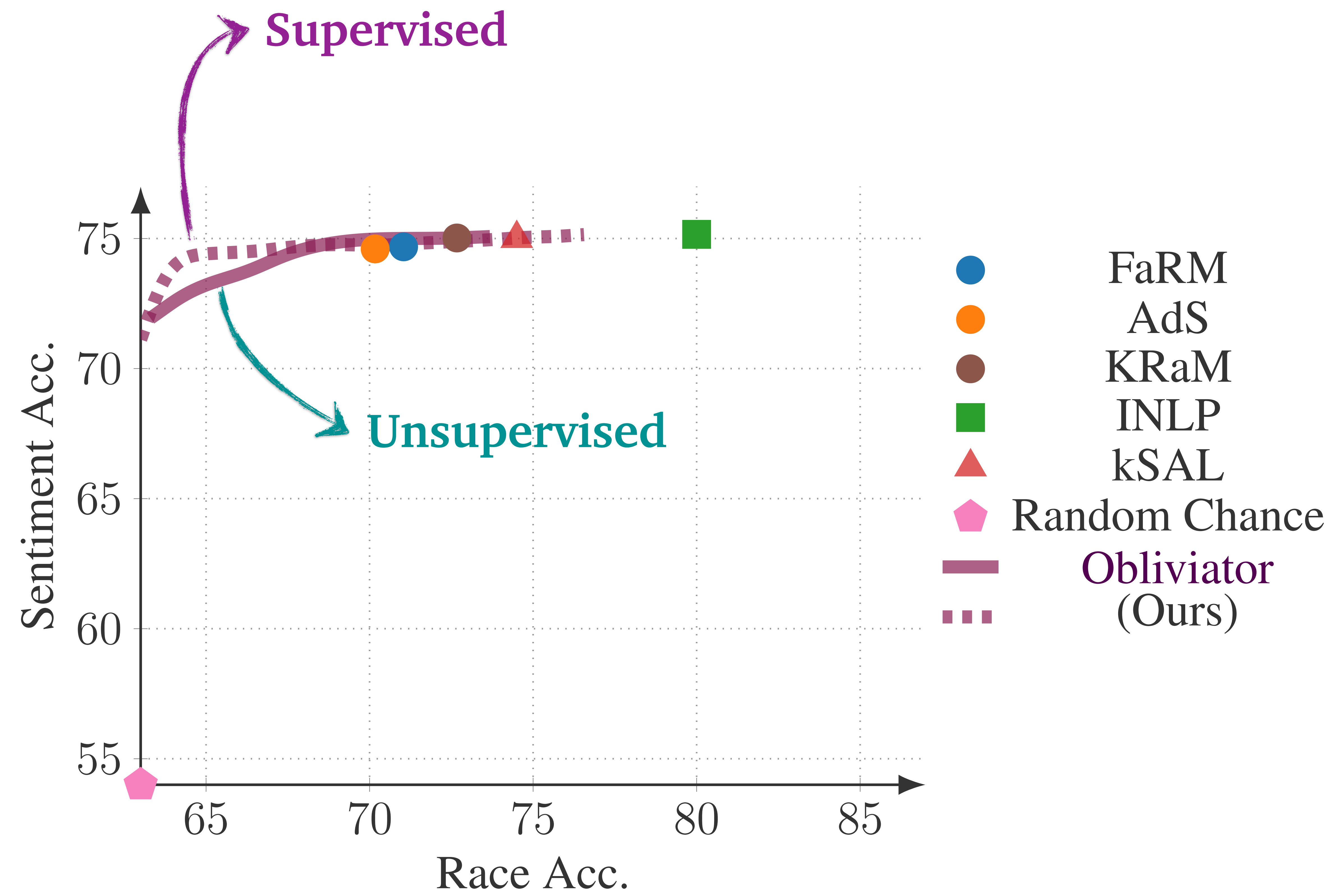
Representation: **Frozen**

PLM : **BERT**

Obliviator Reveals : Effect of **Target Task Labels** on Erasure



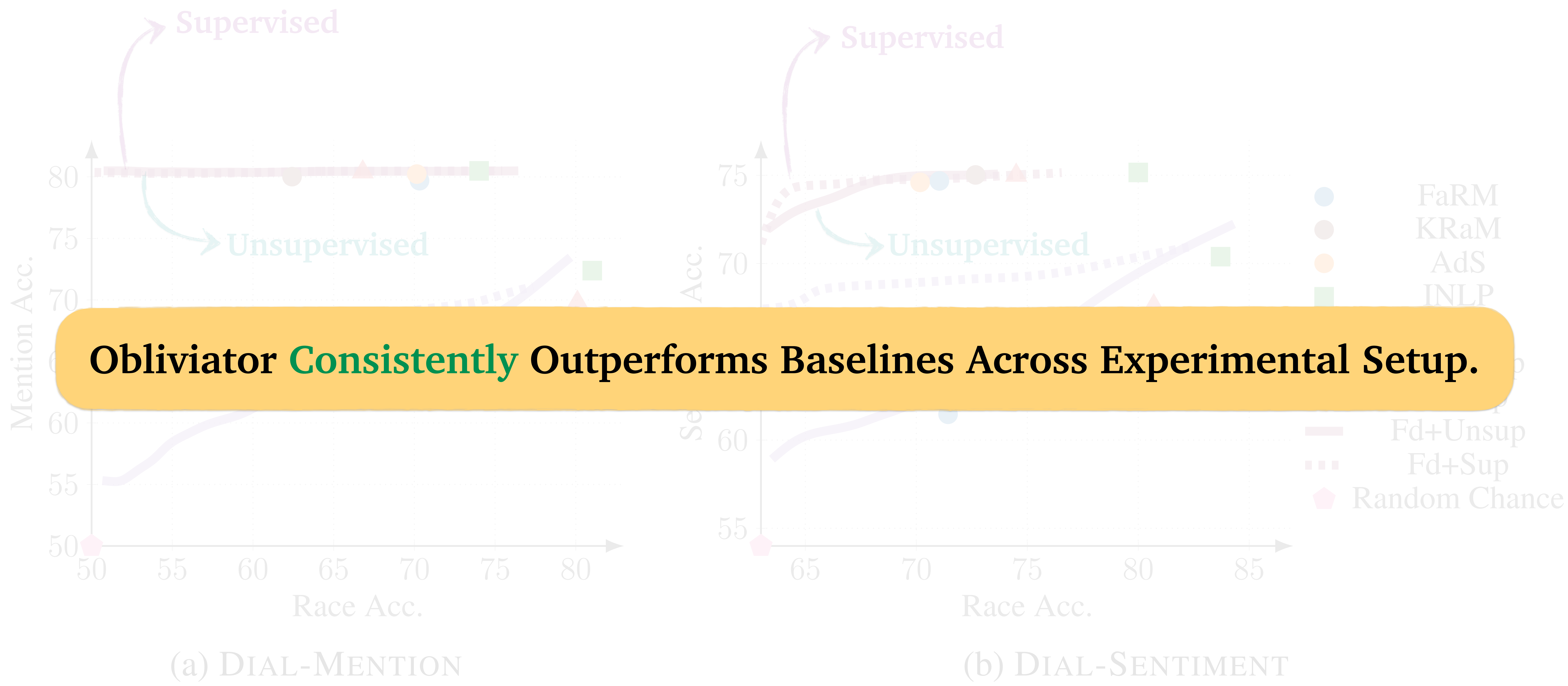
(a) DIAL-MENTION



(b) DIAL-SENTIMENT

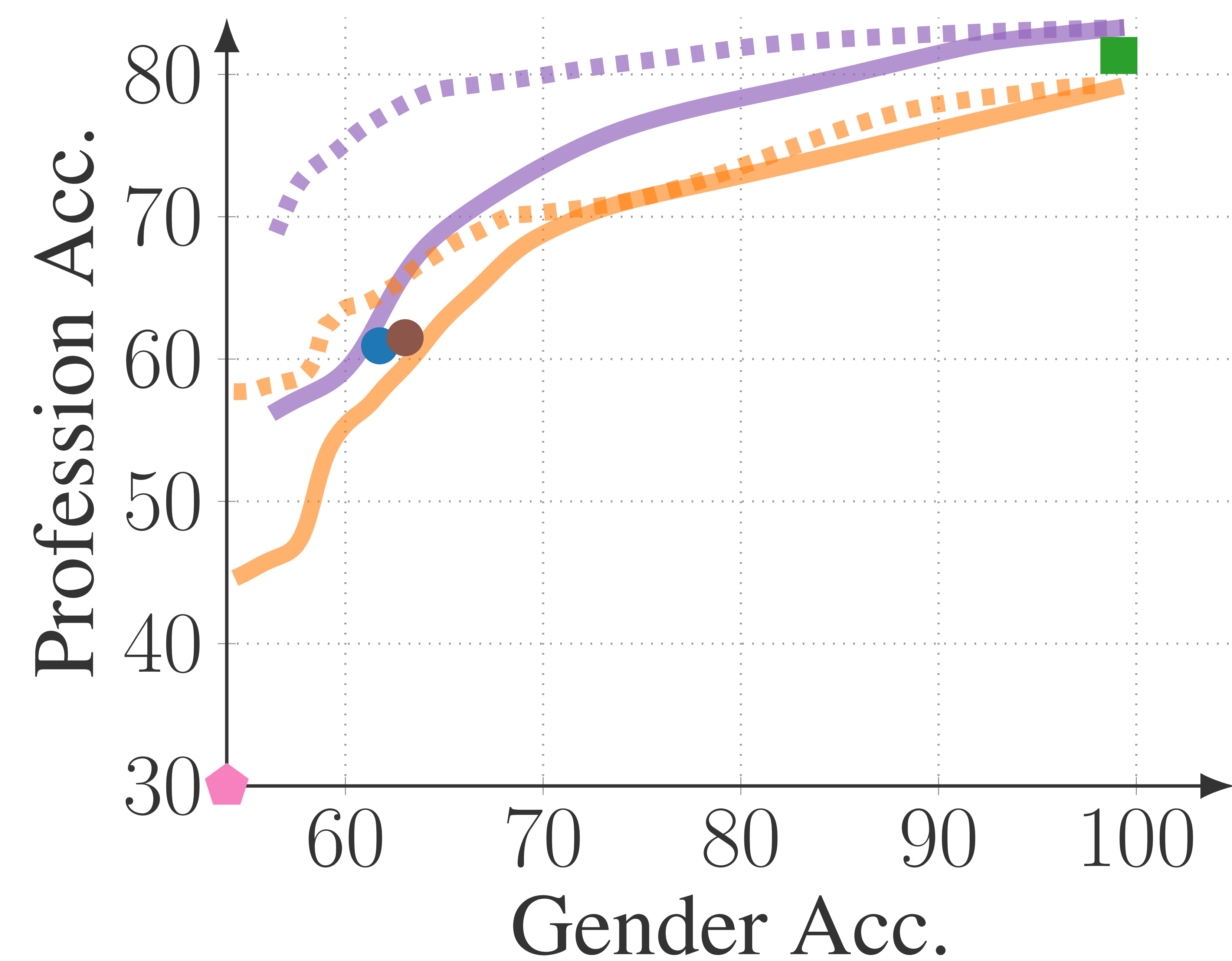
Representation: **Finetuned** PLM : **BERT**

Obliviator Reveals : Effect of **Target Task Labels** on Erasure

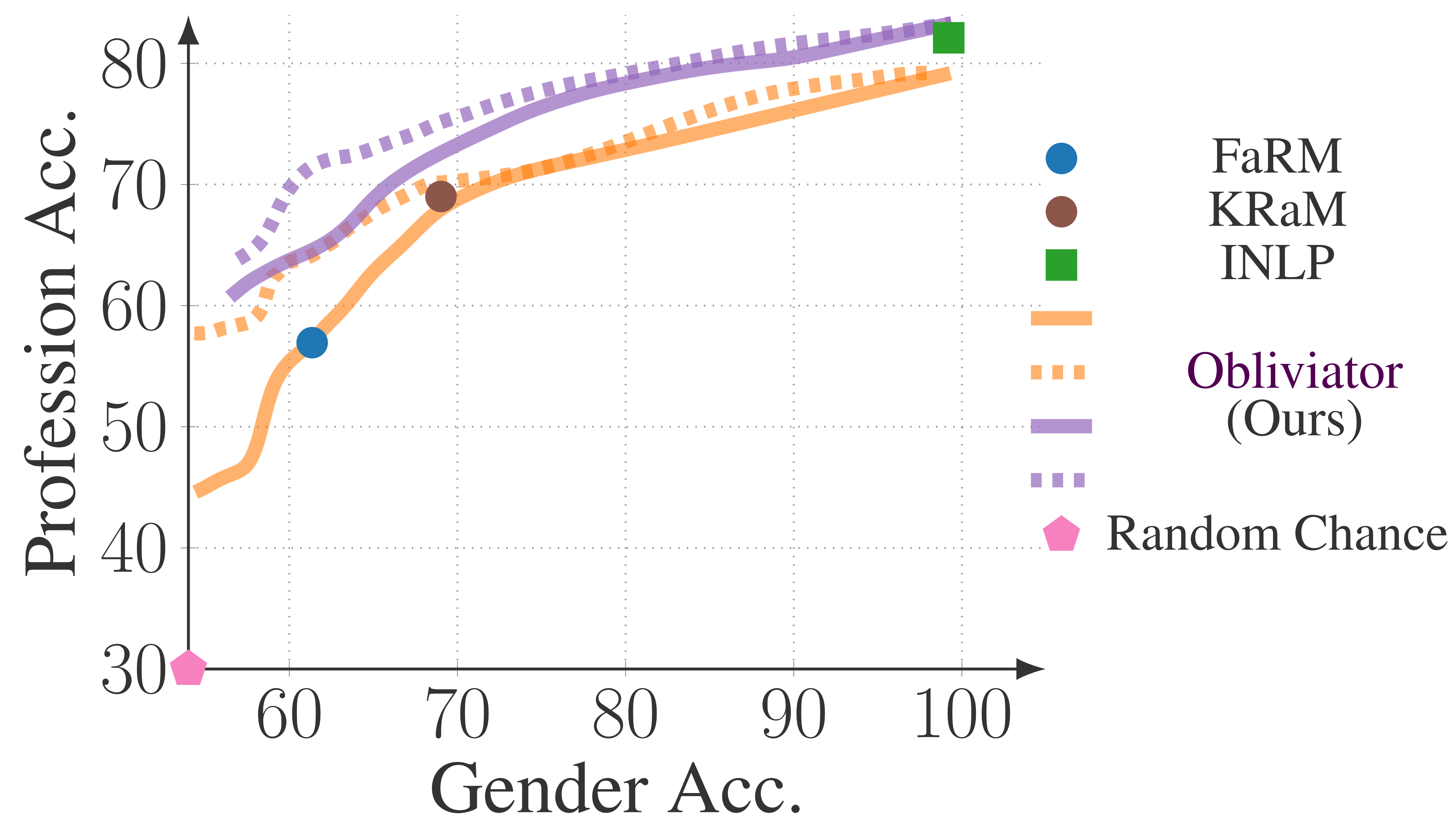


Obliviator Reveals : Erasure Uncovers Representation Structure

Obliviator Reveals : Erasure Uncovers Representation Structure

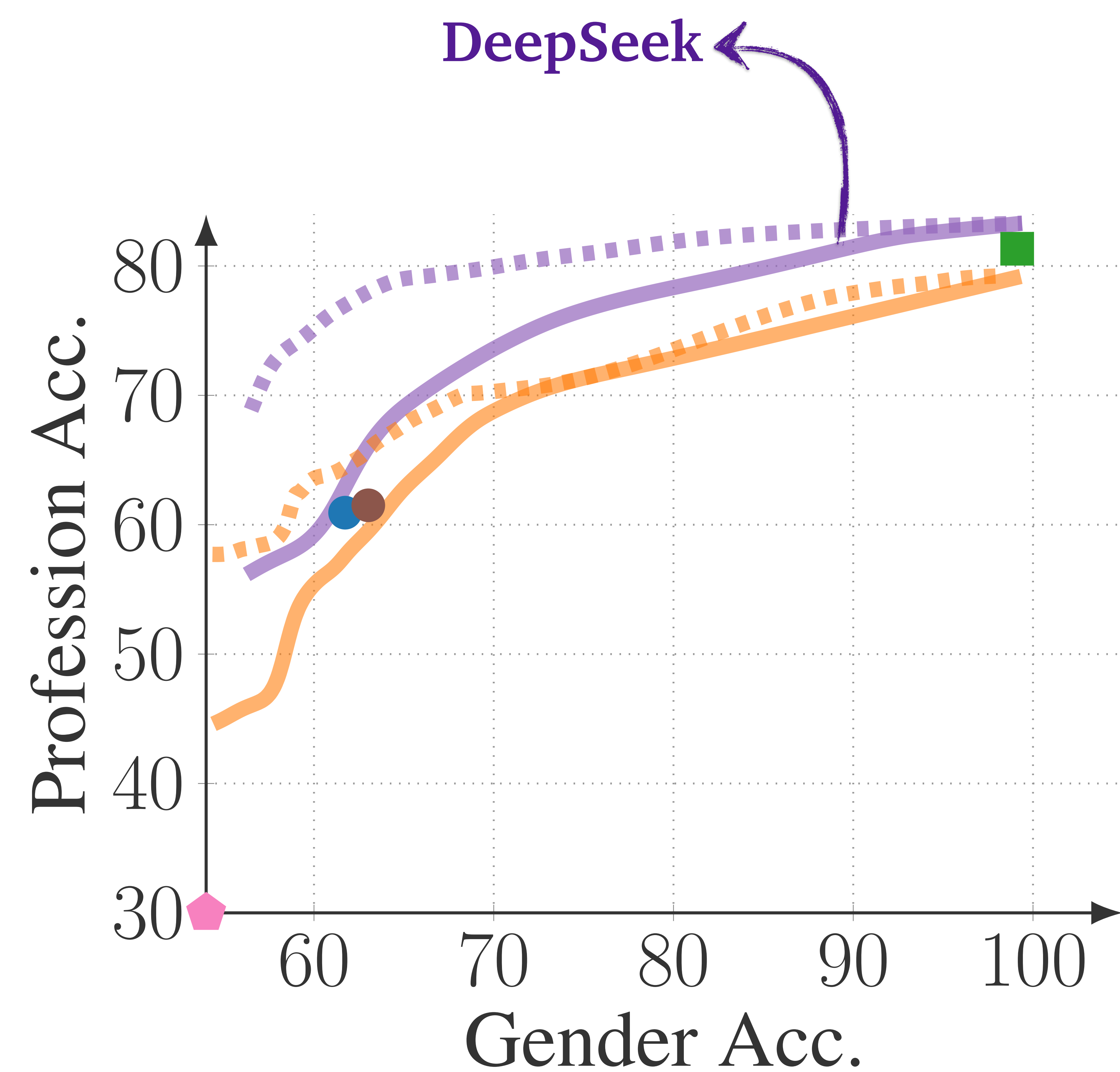


(a) DeepSeek Representations

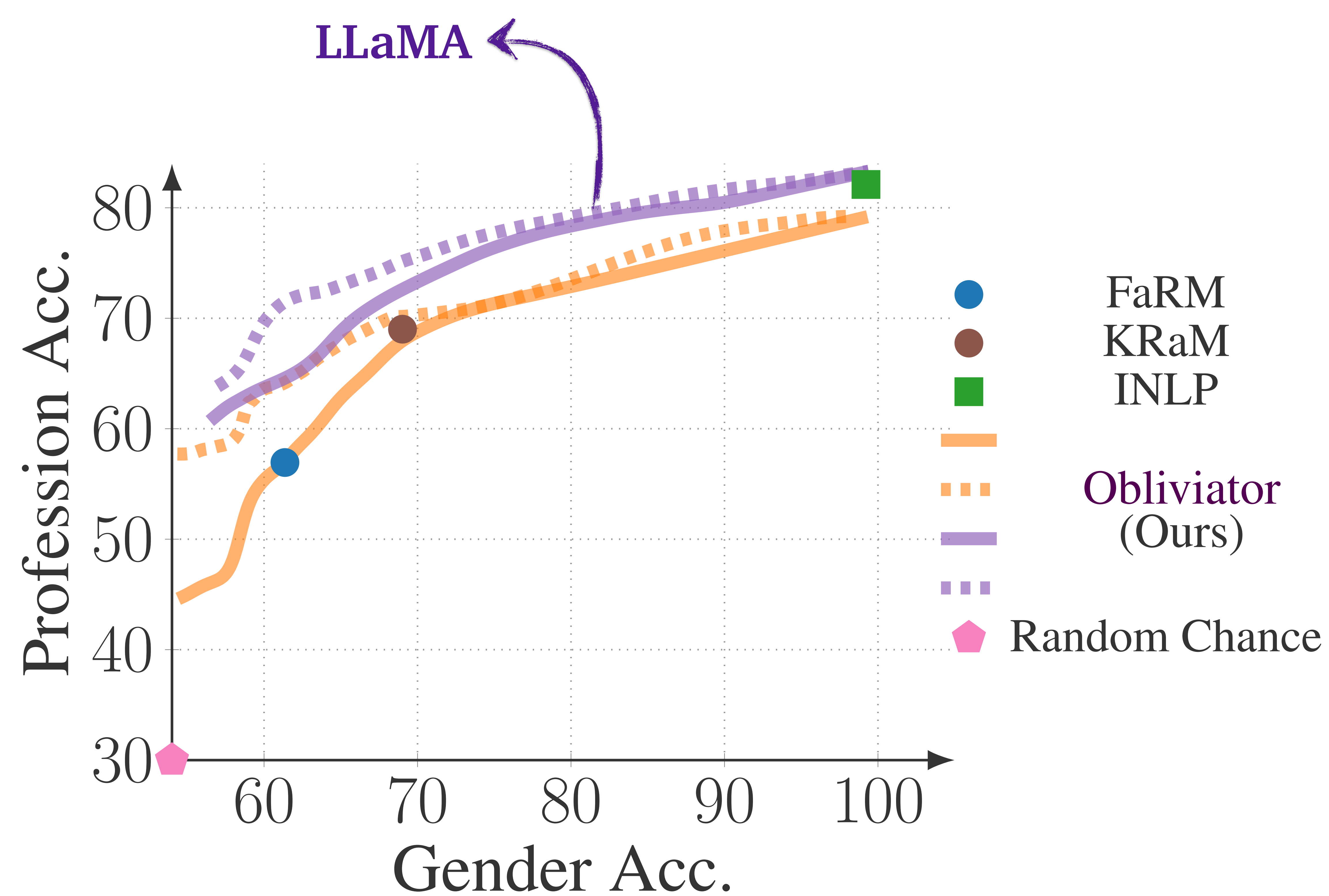


(b) LLaMA Representations

Obliviator Reveals : Erasure Uncovers Representation Structure

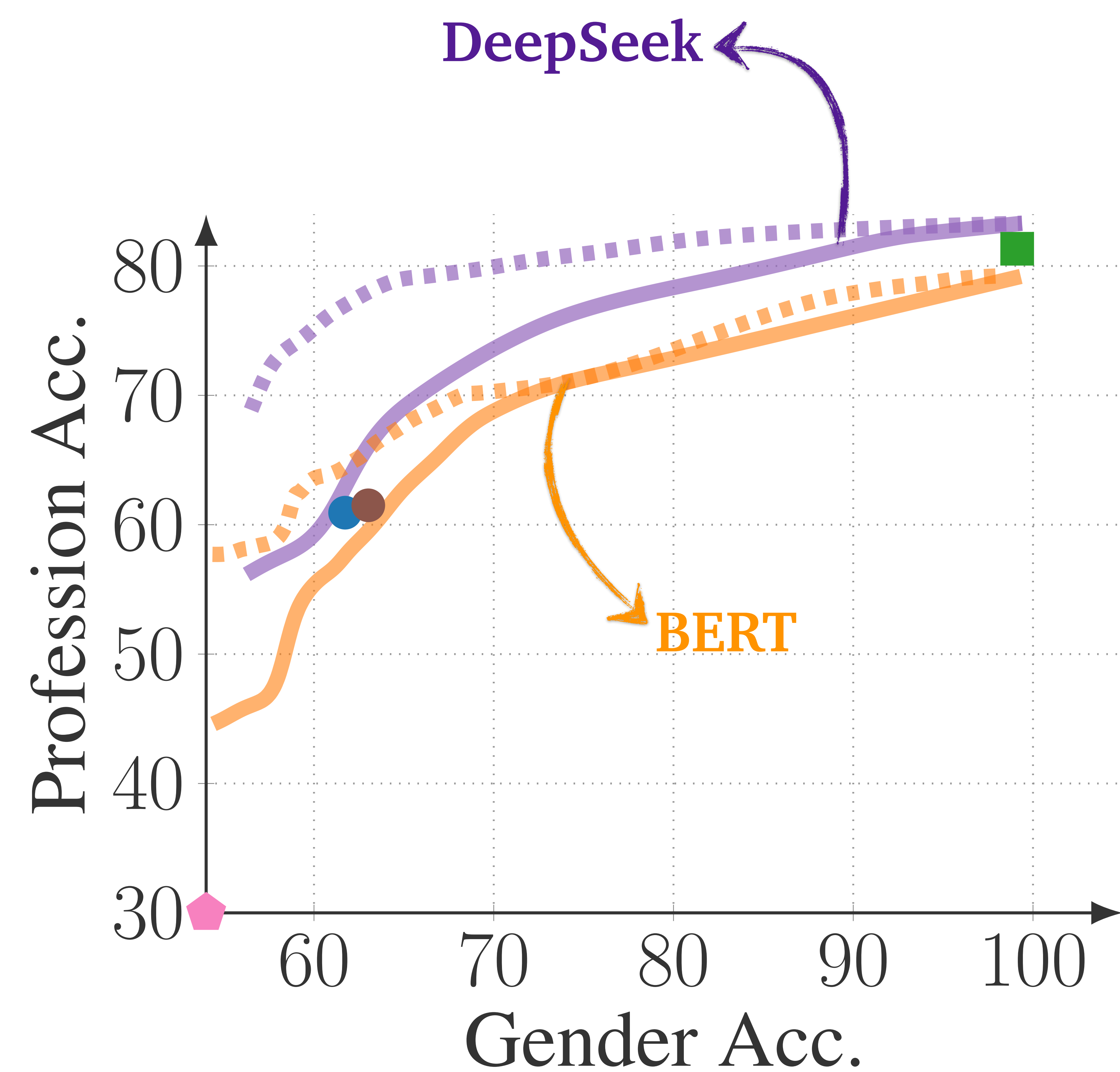


(a) DeepSeek Representations

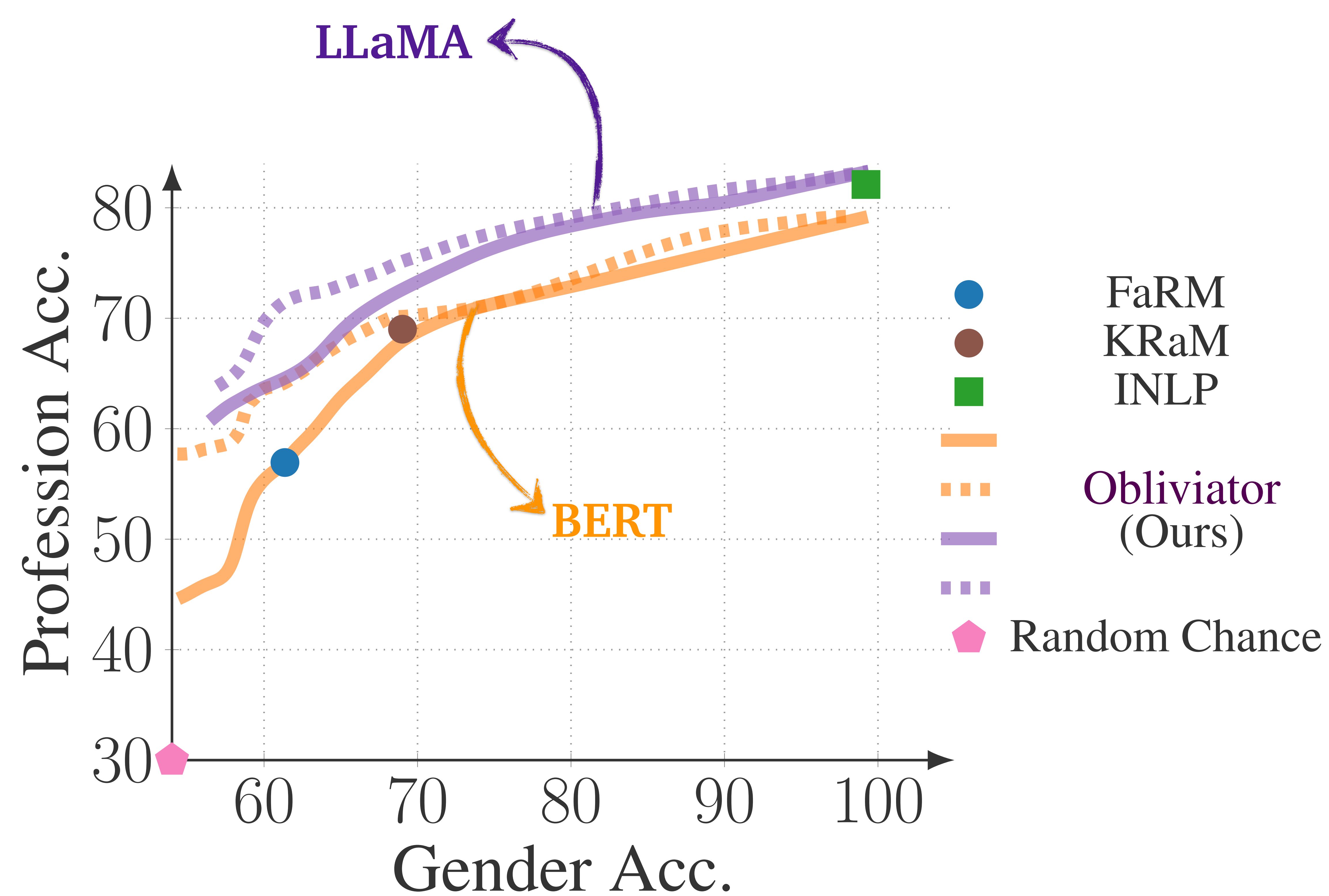


(b) LLaMA Representations

Obliviator Reveals : Erasure Uncovers Representation Structure

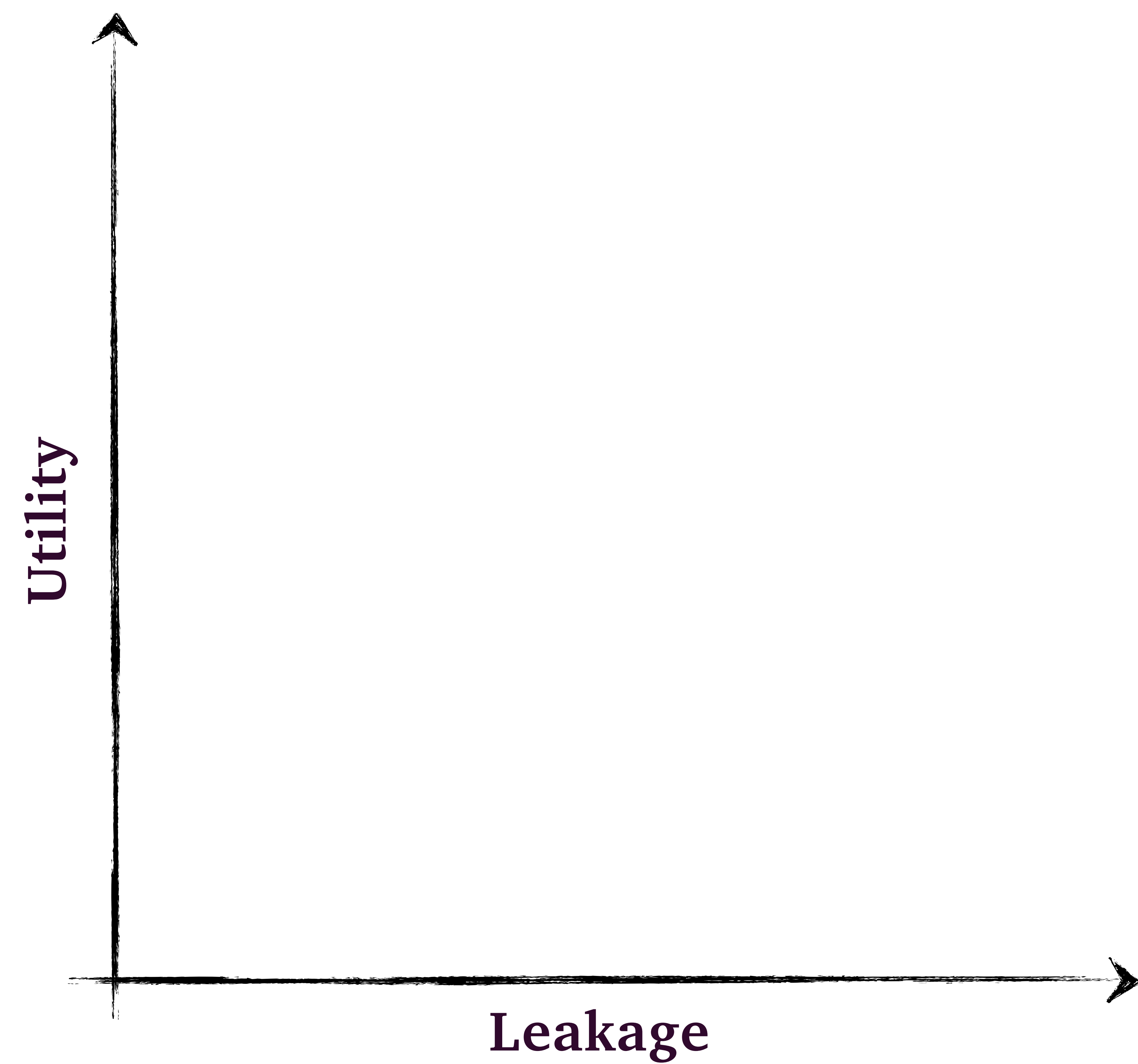


(a) DeepSeek Representations

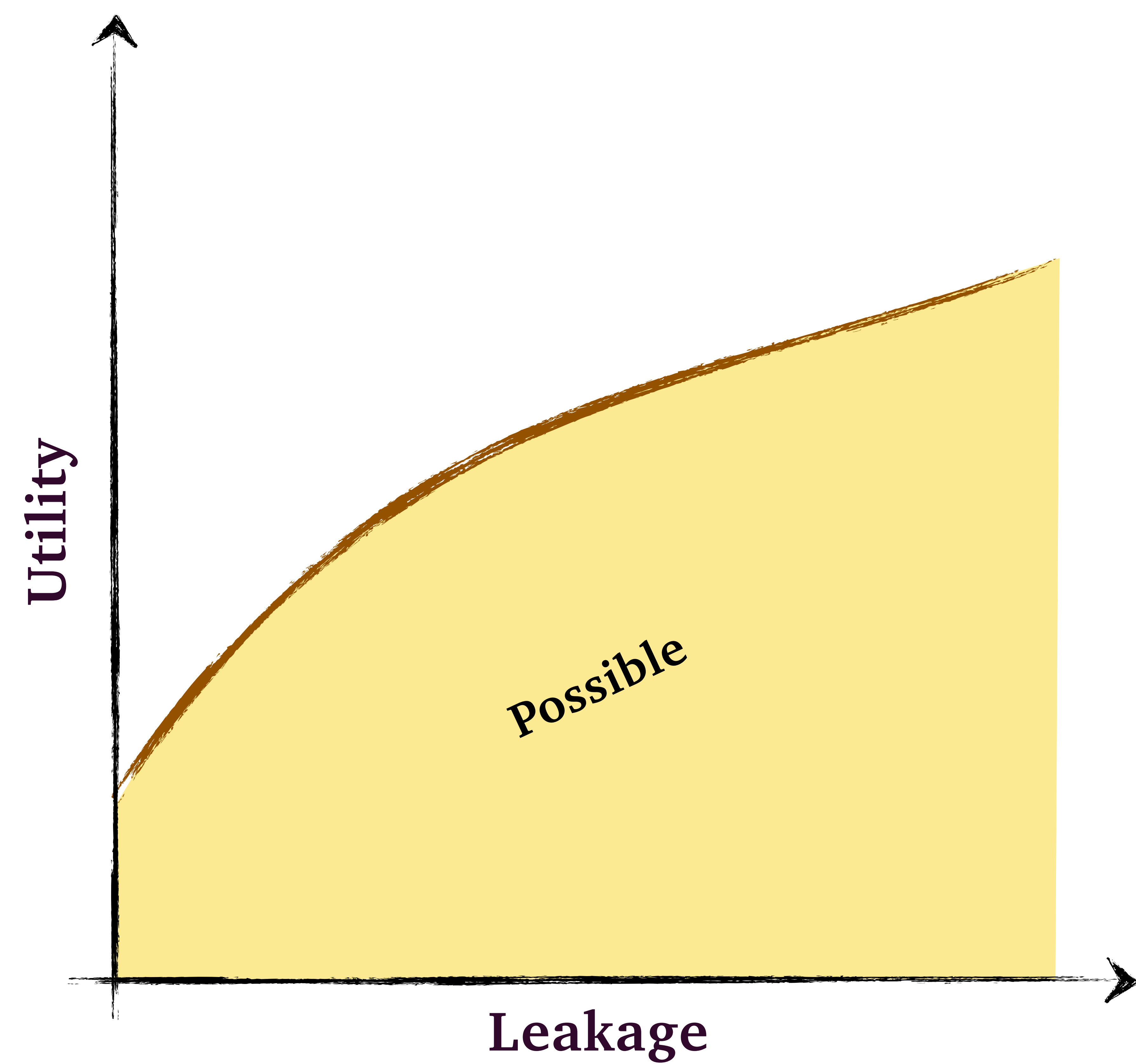


(b) LLaMA Representations

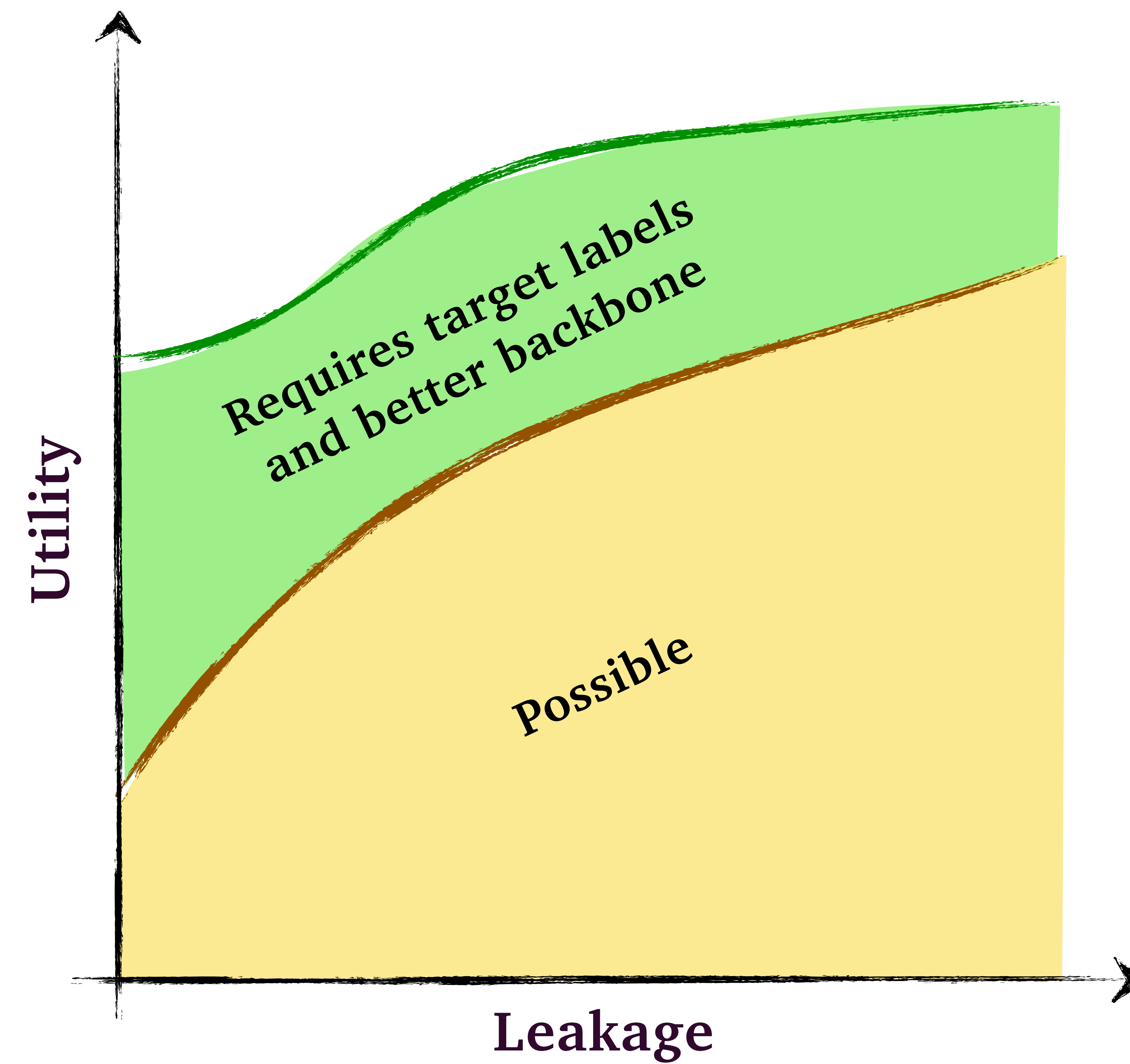
Obliviator Reveals : Utility-Erasure Trade-off



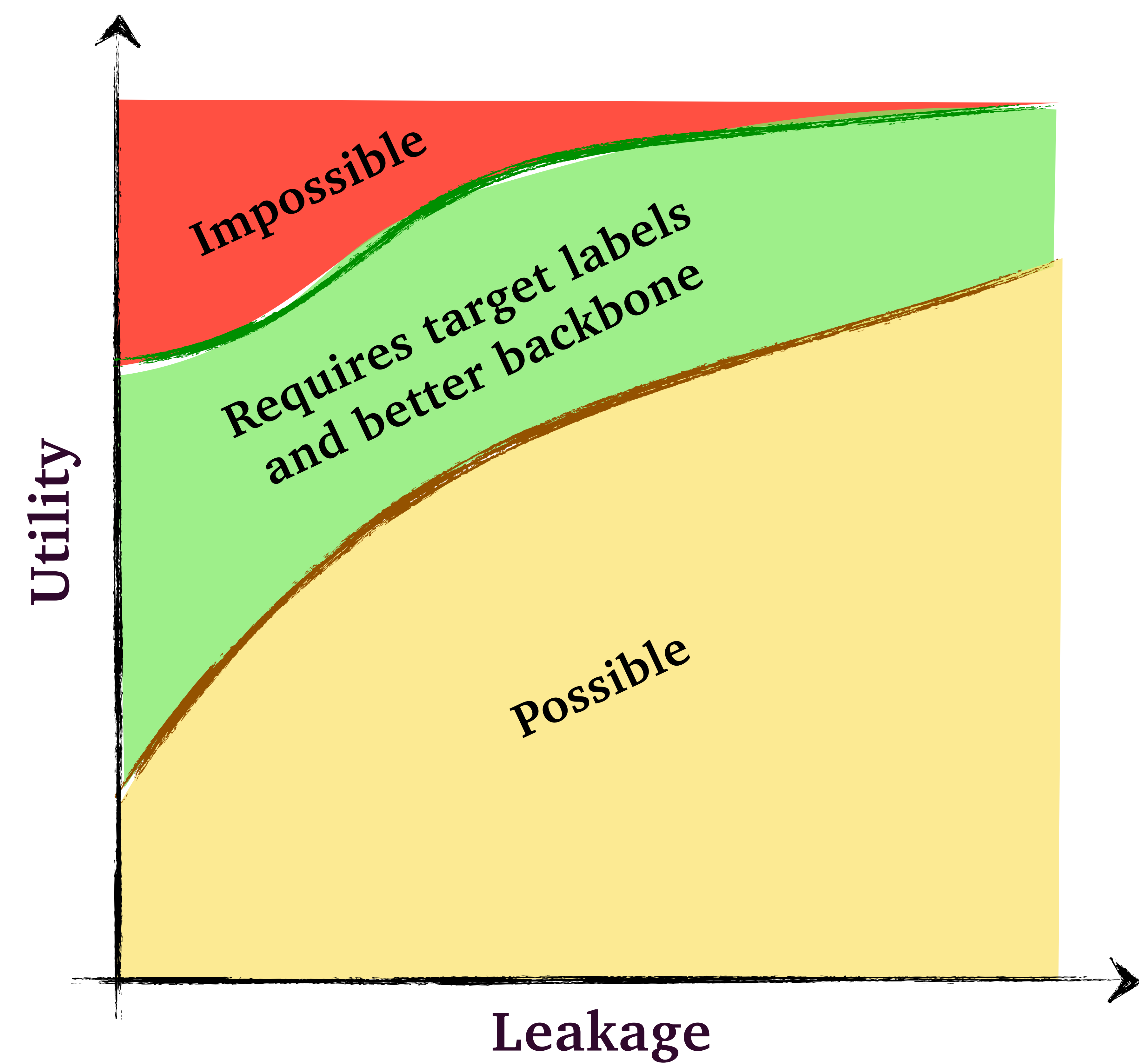
Obliviator Reveals : Utility-Erasure Trade-off



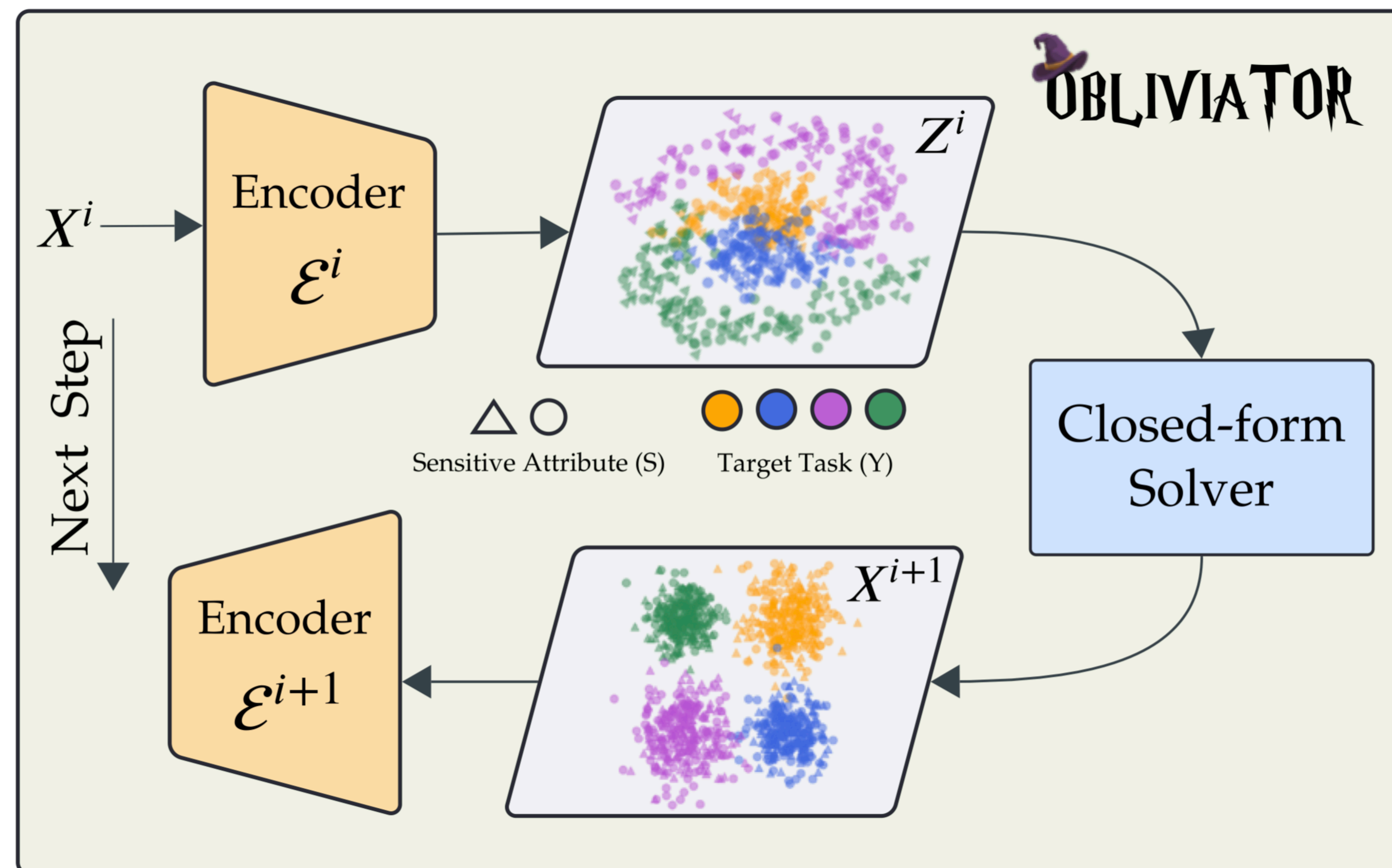
Obliviator Reveals : **Utility-Erasure Trade-off**



Obliviator Reveals : Utility-Erasure Trade-off



Obliviator At a Glance



- **SOTA** Utility-Erasure Trade-off
- Achieves **Nonlinear** Guardedness
- Computationally **Efficient**
- **Fine-Control** over Erasure

