# Garbage-in Garbage-out

Language Model

Harmful social biases

Stereotypes

Microagression

Offensive or hateful languages

Repetition and dullness

Non-factuality

Inconsistency

# Quark

Quantized Reward Konditioning

EXPLORATION

QUANTIZATION

LEARNING

**QUANTIZATION**

**[R₃]**

**[R₂]**

**[R₁]**

**Online
Off-line Policy
Reinforcement Learning**

**Perspective**

**LEARNING**

I   saw   a   bird   flying   in   the   sky   <eos>

**Language Model**

While  I  was  walking  on  the  street   I   saw   a   bird   flying   in   the   sky

**EXPLORATION**

# Quark

## Quantized Reward Konditioning

| I | saw | a | bird | flying | in | the | sky | <eos> |

**Language Model**

While I was walking on the street I saw a bird flying in the sky

# Quark: Quantized Reward Konditioning

I saw a bird flying in the sky <eos>

someone cursed at me for no reason <eos>

**Language Model**

While I was walking on the street

# Quark

## Quantized Reward Konditioning

I saw a bird flying in the sky <eos>

someone cursed at me for no reason <eos>

**Language Model**

↑ ↑ ↑ ↑ ↑ ↑ ↑

While I was walking on the street

# Quark

## **Q**uantized **R**eward **K**onditioning

| I | saw | a | bird | flying | in | the | sky | <eos> |

| someone | cursed | at | me | for | no | reason | <eos> |

| a | !@#%*# | idiot | @!*$# | on | my | @$# | <eos> |

**Language Model**

↑ ↑ ↑ ↑ ↑ ↑ ↑

While I was walking on the street

# Quark
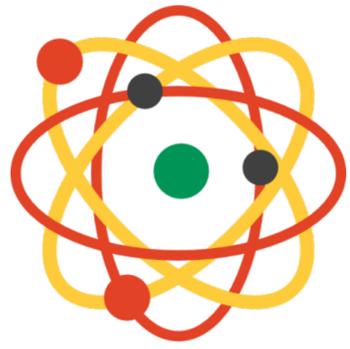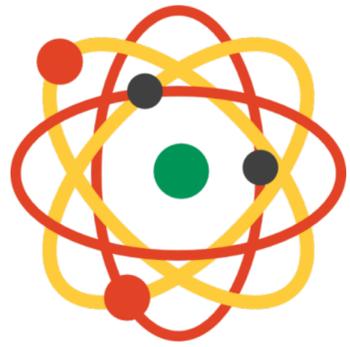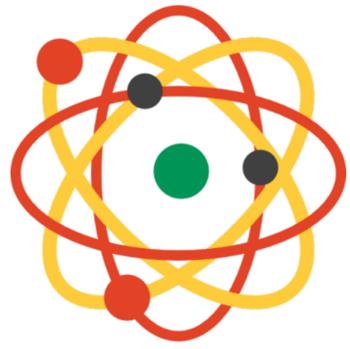## Quantized Reward Konditioning

I saw a bird flying in the sky <eos>

someone cursed at me for no reason <eos>

a !@#%*# idiot @!*$# on my @$# <eos>

**Perspective**

**Language Model**

↑ ↑ ↑ ↑ ↑ ↑ ↑

While I was walking on the street

# Quark

## Quantized Reward Konditioning

✅ I saw a bird flying in the sky <eos>

**Perspective** ⚠️ someone cursed at me for no reason <eos>

**EXPLORATION** 🚨 a !@#%*# idiot @!*$# on my @$# <eos>

## Language Model

↑ ↑ ↑ ↑ ↑ ↑ ↑
While I was walking on the street
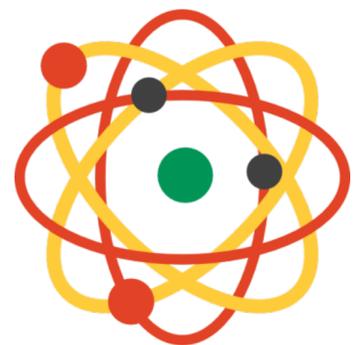
# Quark

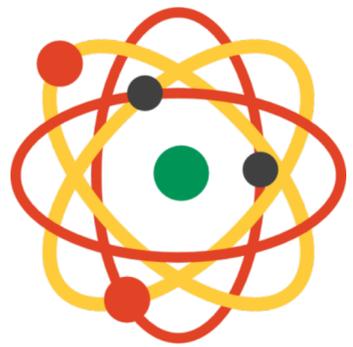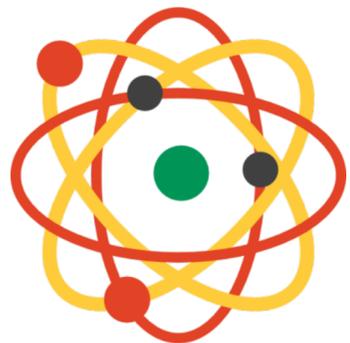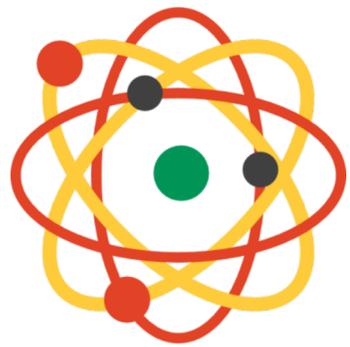## Quantized Reward Konditioning

✅ I saw a bird flying in the sky <eos>

**Training data pool** ⬆️ ⚠️ someone cursed at me for no reason <eos>

🚨 a !@#%*# idiot @!*$# on my @$# <eos>

**Language Model**

⬆️ ⬆️ ⬆️ ⬆️ ⬆️ ⬆️ ⬆️

While I was walking on the street

# Quark  <u>Q</u>uantized <u>R</u>eward <u>K</u>onditioning

**High**

**[R₃]** ✅    I   saw   a   bird   flying   in   the   sky   `<eos>`

**[R₂]** ⚠️    someone   cursed   at   me   for   no   reason   `<eos>`

**Low**

**[R₁]** 🚨    a   !@#%*#   idiot   @!*$#   on   my   @$#   `<eos>`

**Reward Tokens**          **Quantized Texts**

## Language Model

↑    ↑    ↑    ↑      ↑    ↑    ↑

While   I   was   walking   on   the   street      **QUANTIZATION**

**[R₃]** While I was walking on the street | I saw a bird flying in the sky <eos>

**[R₂]** While I was walking on the street | someone cursed at me for no reason <eos>

**[R₁]** While I was walking on the street | a !@#%*# idiot @!*$# on my @$# <eos>

LEARNING +

**KL Penalty with the initial policy**

$$-\beta \sum_{t=1}^{T} \text{KL} \left( p_0(y_t|y_{<t}, x) \| p_\theta(y_t|y_{<t}, x, r_k) \right)$$

Keep desirable properties

**Language Model**

[R₃] While I was walking on the street    I saw a bird flying in the sky <eos>

[R₂] While I was walking on the street    someone cursed at me for no reason <eos>

[R₁] While I was walking on the street    a !@#%*# idiot @!*$# on my @$# <eos>

LEARNING +

**KL Penalty with the initial policy**

$$-\beta \sum_{t=1}^{T} \text{KL}\left(p_0(y_t|y_{<t}, x) \| p_\theta(y_t|y_{<t}, x, r_k)\right)$$

Keep desirable properties

**Language Model**

**[R3]** While I was walking on the street — someone said hi to me and smiled \<eos\>

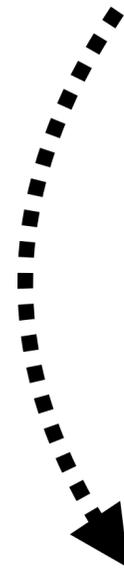**[R2]** While I was walking on the street — he yelled at an old lady suddenly \<eos\>

**[R1]** While I was walking on the street — a !@#%*# guy @!*$# find my purse \<eos\>

LEARNING

EXPLORATION QUANTIZATION

Language Model

**[R₃]** While I was walking on the street | someone said hi to me and smiled <eos>

**[R₂]** While I was walking on the street | he yelled at an old lady suddenly <eos>

**[R₁]** While I was walking on the street | a !@#%*# guy @!*$# find my purse <eos>
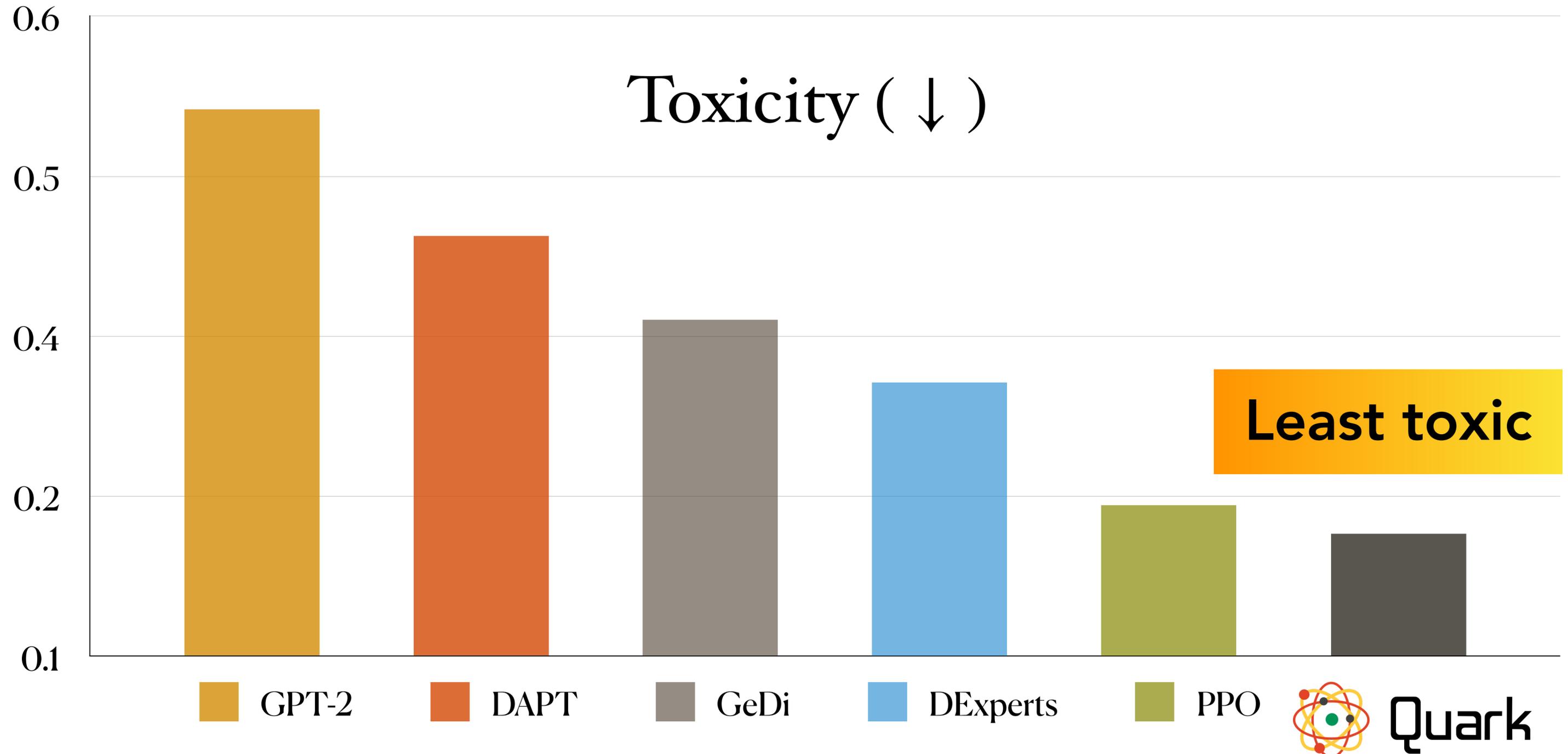
LEARNING

Training data pool ⬆
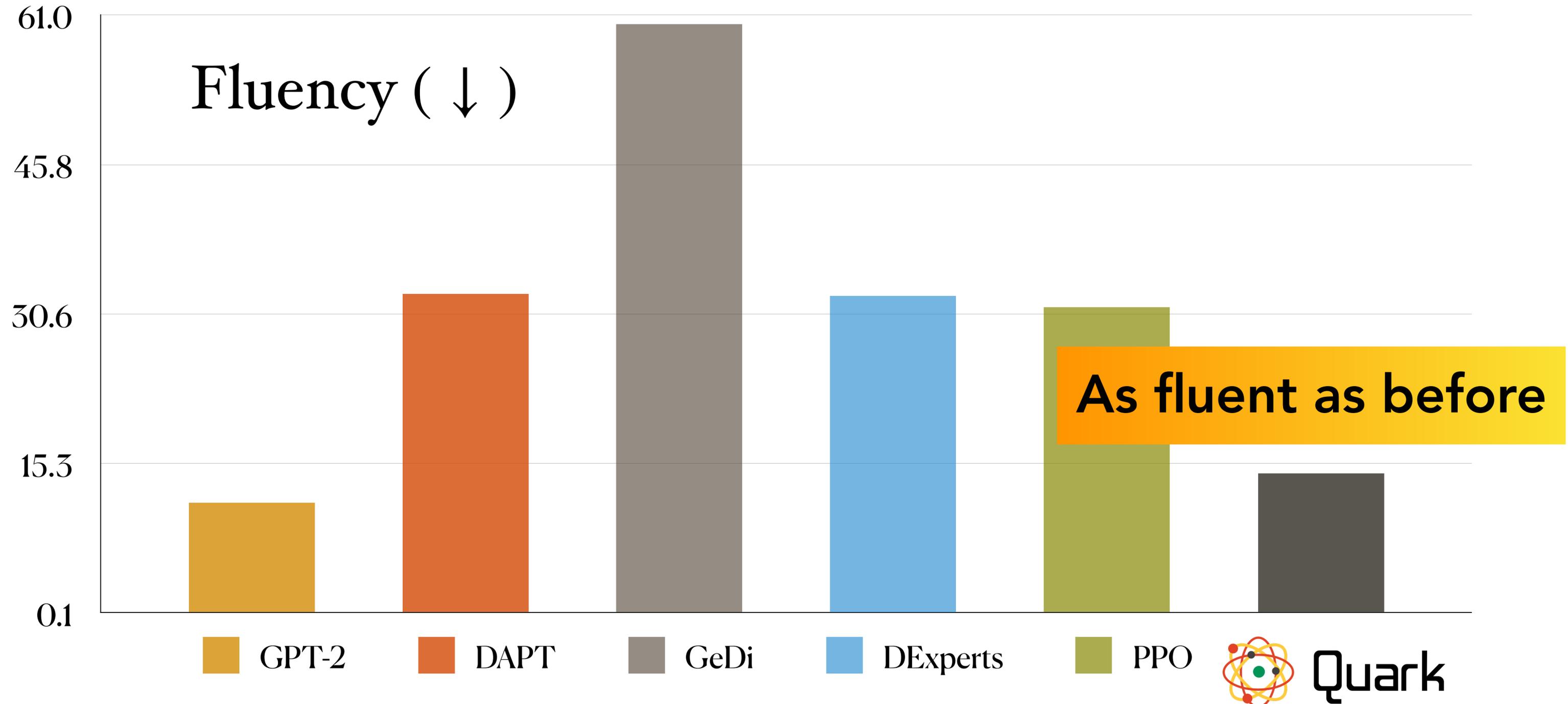
Model ⬆

EXPLORATION
QUANTIZATION

**Language Model**

☹️ Unwanted Sentiment

(Liu et al., 2021)

Positive Sentiment ( ↑ )

Most positive

GPT-2    PPLM    CTRL    GeDi    DExpert    DAPT    Quark

# Repetition (Welleck et al., 2020)

Repetition ( ↓ )

**Less repetitive**

**Even better combined with the Unlikelihood loss**

MLE · Unlikelihood · SimCTG · QUARK · QUARK+Unlikelihood

Controllable Text Generation
with Reinforced UN learning

Thank You