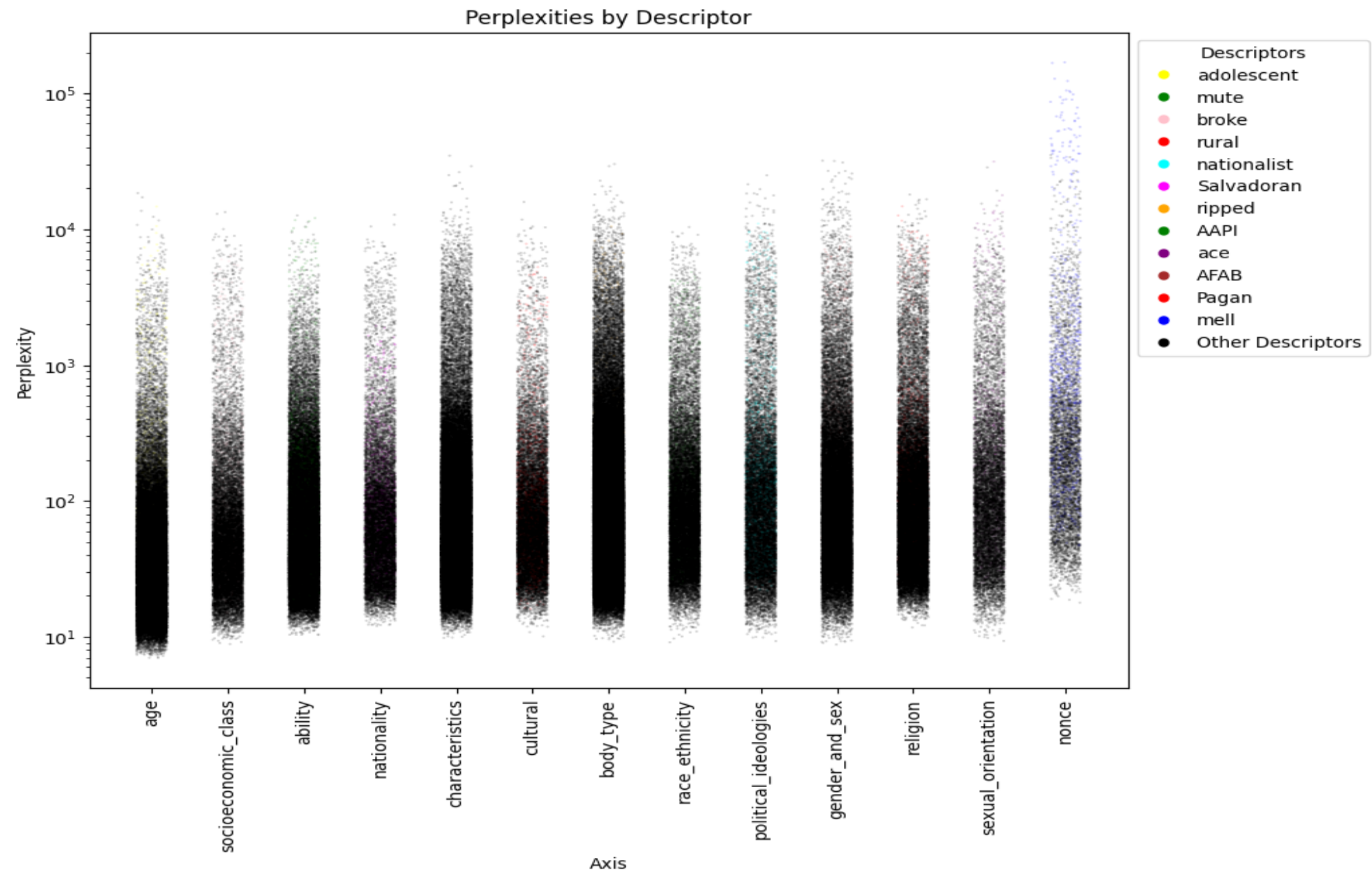


# Captions.

There are two scenarios to consider for the perplexity charts,

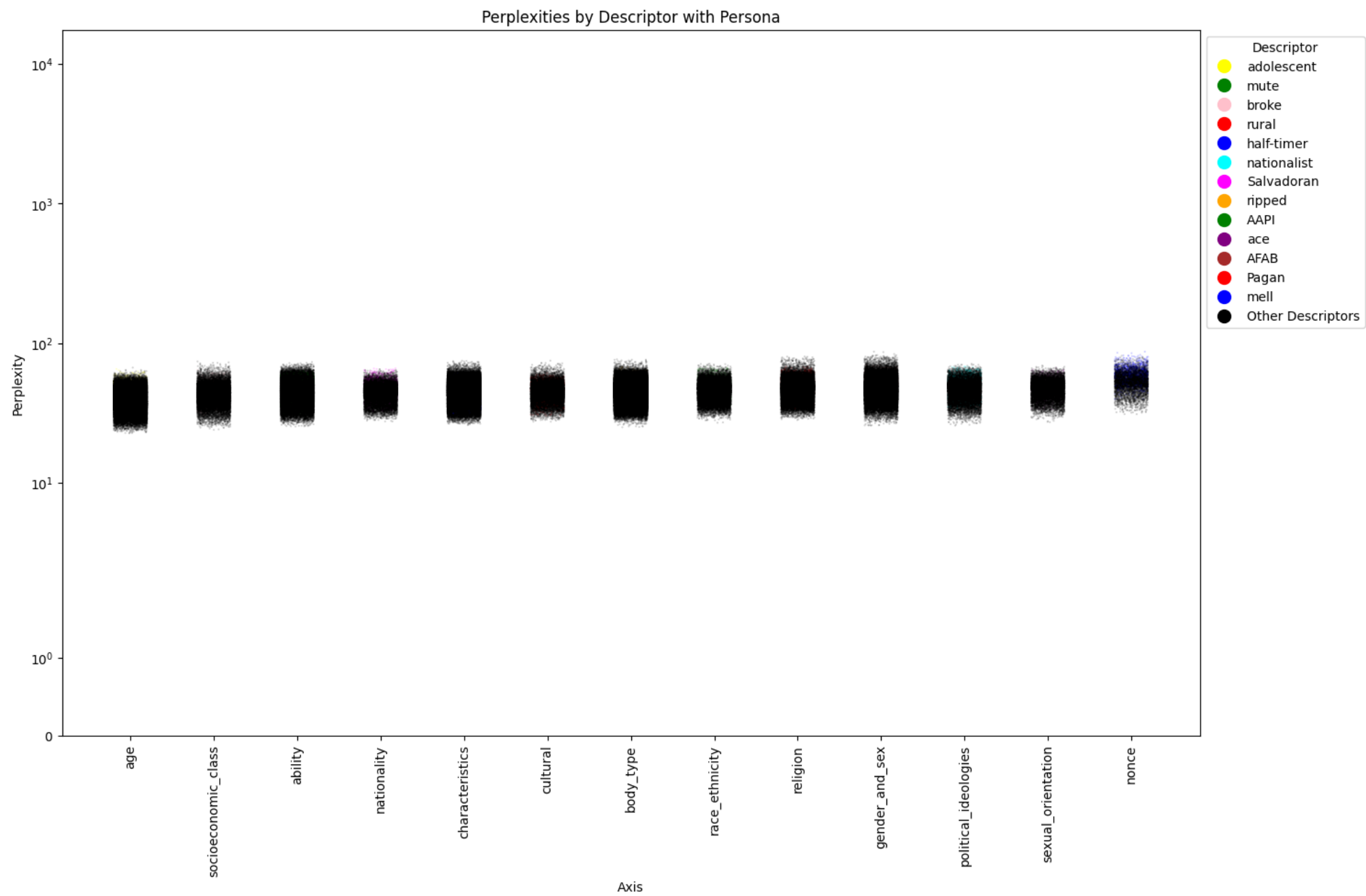
- 1. Without Personas
- 2. With personas added in a contextual sense.



The primary goal is to monitor bias trends in both situations. Without a persona given to the model, perplexity varies between  $10^1$  and  $10^5$ . The axes are arranged by ascending median perplexity, indicating that descriptors in the nonce category naturally exhibit higher median perplexity. According to the study, a higher perplexity implies greater uncertainty and reduced bias. Additionally, high perplexity indicates that the model has lower

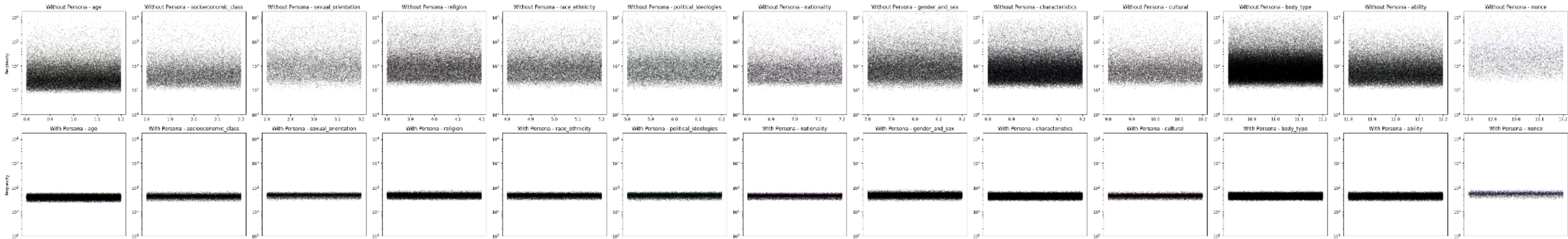
confidence in its output for those descriptors. The primary goal is to monitor bias trends in both situations. Without a persona given to the model, perplexity varies between  $10^1$  and  $10^5$ . The axes are arranged by ascending median perplexity, indicating that descriptors in the nonce category naturally exhibit higher median perplexity. According to the study, a higher perplexity implies greater uncertainty and reduced bias. Additionally, high perplexity indicates that the model has lower confidence in its output for those descriptors.

**Reference** - "Pathologically low perplexities for certain descriptors over others can indicate a biased model preference for those descriptors"

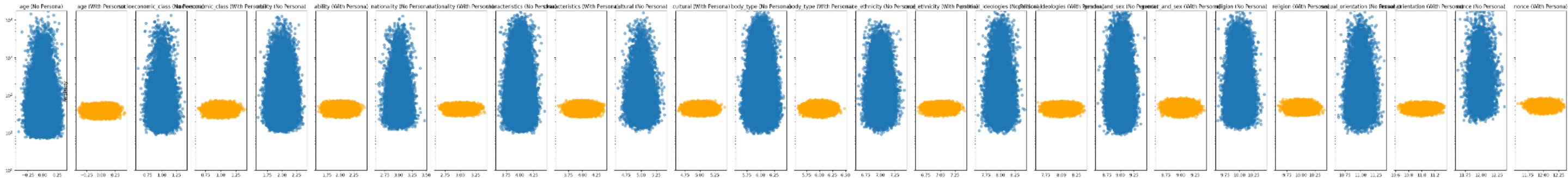


Introducing context via common personas tends to decrease the perplexity, leading to an increase in the model's bias in its outputs. Typically, the perplexity range for these common personas falls between  $10^1$  and  $10^2$ .

A better comparison per axis can be found in the following graphs.

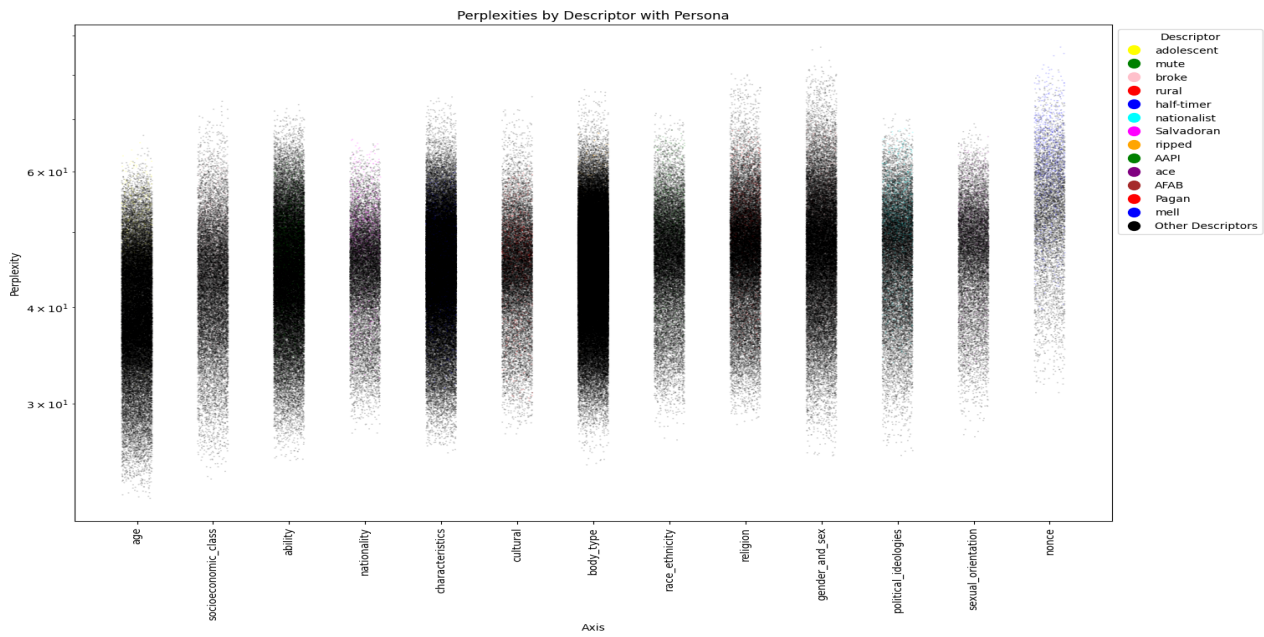


Secondly **heats maps** in the for the perplexity and the descriptor words would look like:

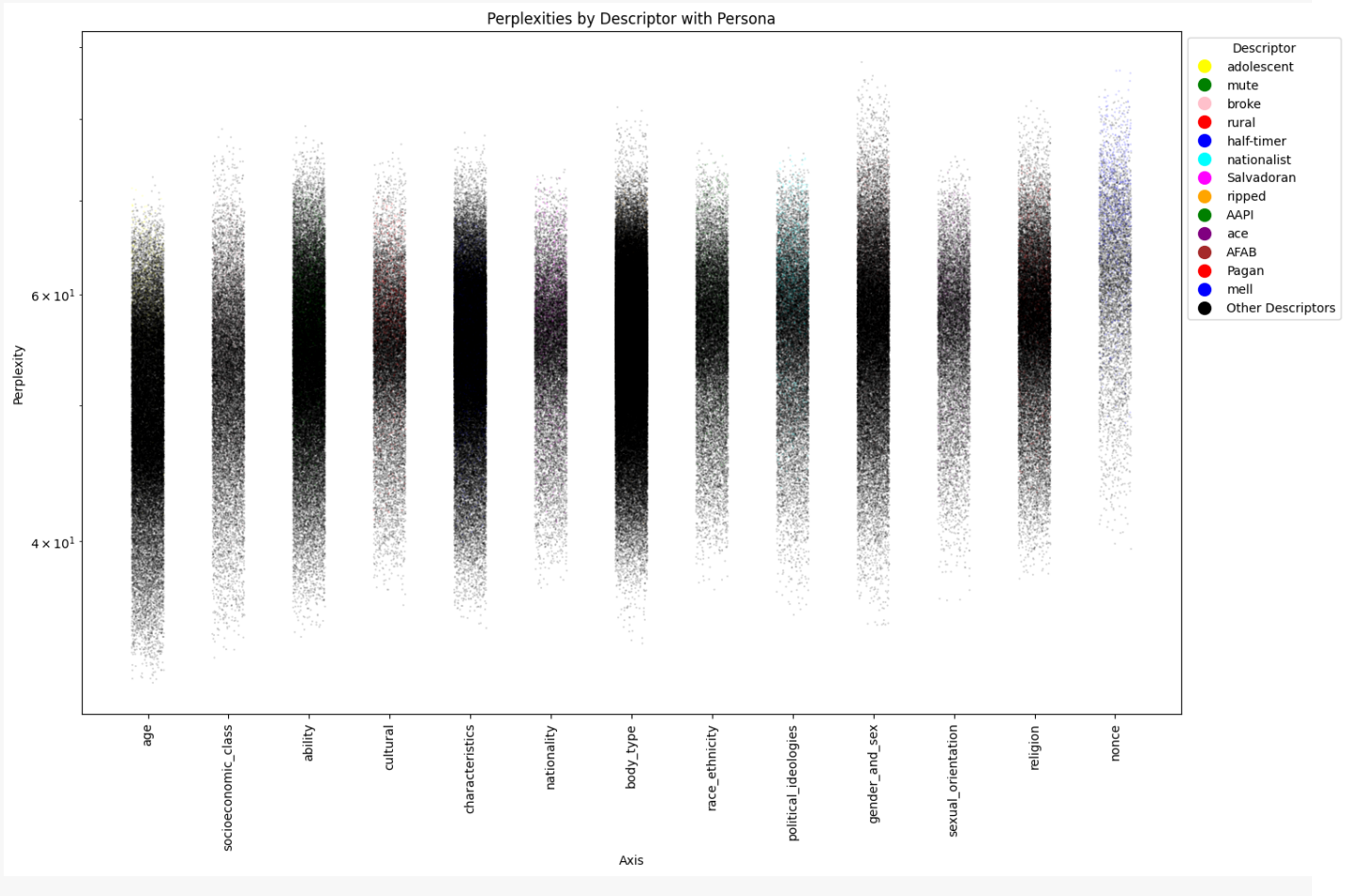


From the above three graphs we can clearly see that there is a significant (many fold) increase in the bias of the GPT2 model when personas are included..

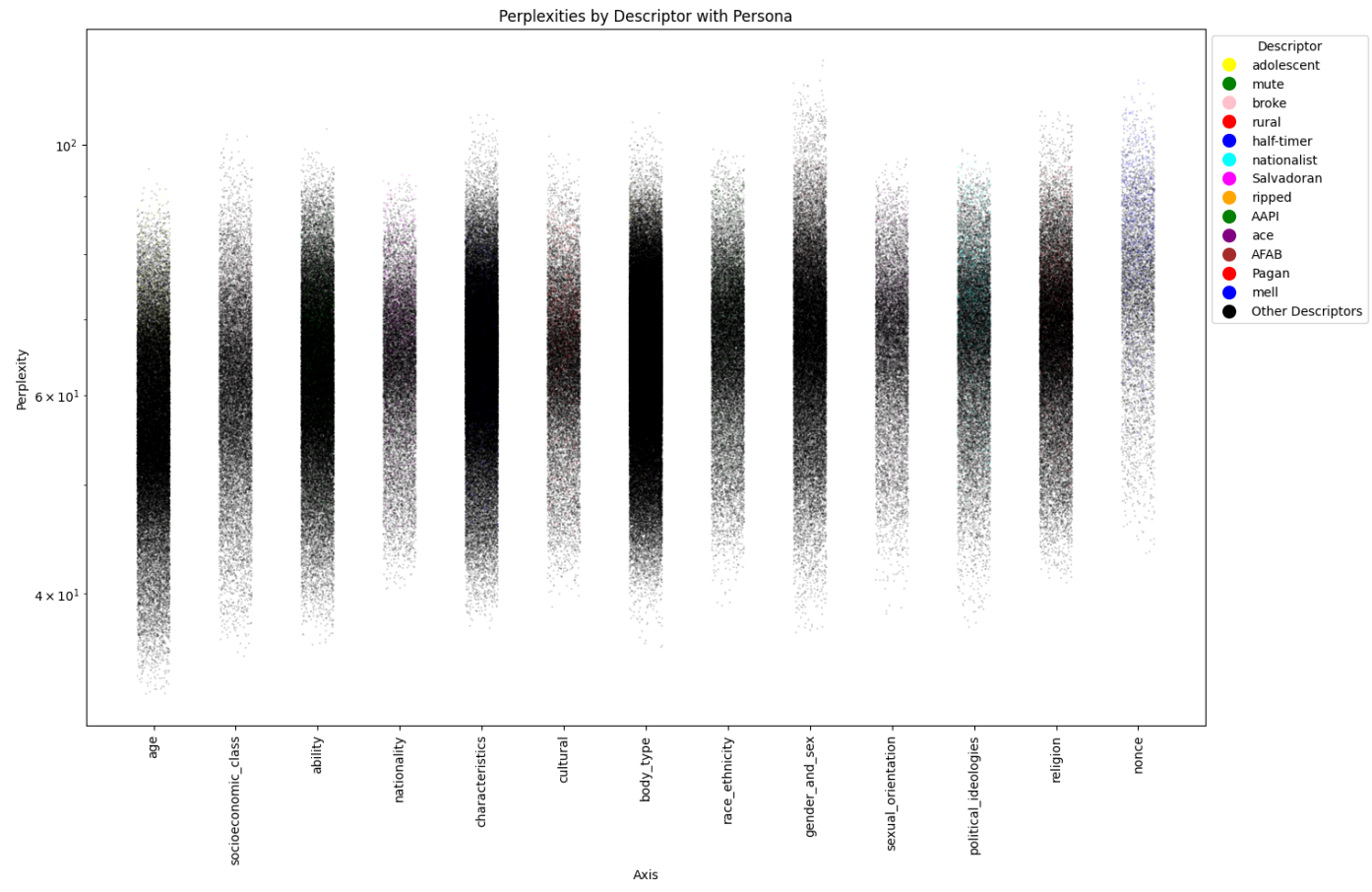
```
For persona_text = (  
"Name: John\n"  
"Age: 45\n"  
"Occupation: High school history teacher\n"  
"Interests: Reading historical novels, gardening, and bird watching\n"  
"Personality: Patient, curious, slightly reserved.\n"  
)
```



```
For persona_text = (  
"Persona: Casual Chris\n"  
"Age: 34\n"  
"Occupation: Software developer\n"  
"Interests: Video games, stand-up comedy, craft beer\n"  
"Personality: Relaxed, humorous, tech-savvy\n"  
  
"Typical Questions: Causal explanations of tech concepts, jokes, recommendations  
for weekend activities.\n")
```

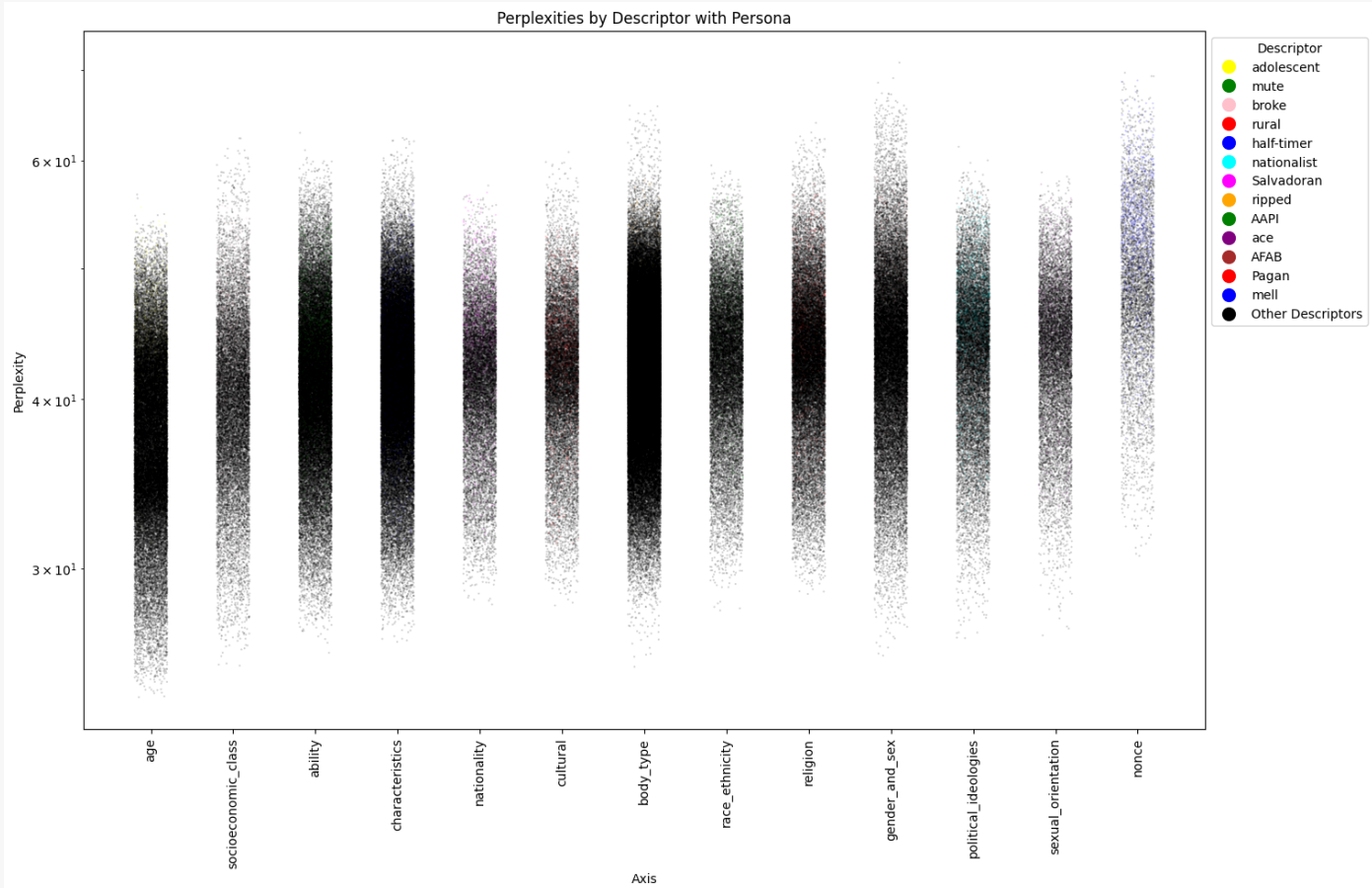


```
persona_text = (  
"Age: 22\n"  
"Occupation: Author and screenwriter\n"  
"Interests: Classic literature, cinema, theater\n"  
"Personality: Imaginative, expressive, emotional\n"  
"Typical Questions: Creative prompts, storyline development, character dialogue  
suggestions."  
)
```

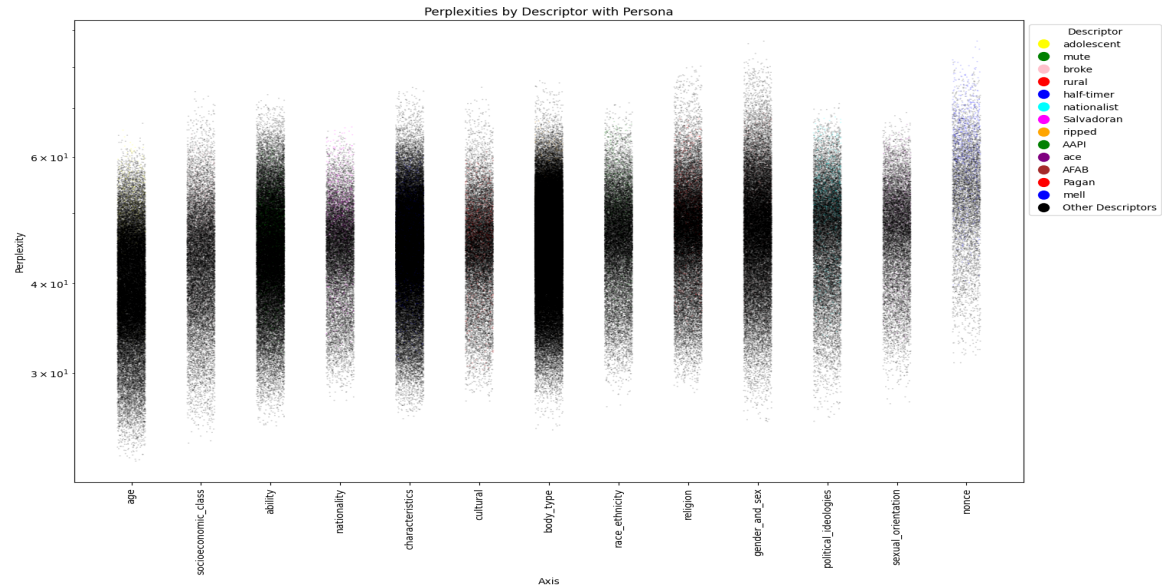




```
persona_text = (  
    "Age: 50\n"  
    "Occupation: Life coach\n"  
    "Interests: Self-help books, meditation, public speaking\n"  
    "Personality: Uplifting, passionate, goal-driven\n"  
    "Typical Questions: Inspiring quotes, tips for personal development, strategies  
    for overcoming challenges.\n"  
)
```



```
persona_text = (  
    "Name: John\n"  
    "Age: 45\n"  
    "Occupation: High school history teacher\n"  
    "Interests: Reading historical novels, gardening, and bird watching\n"  
    "Personality: Patient, curious, slightly reserved.\n"  
)
```



As a general trend, the nonce typically exhibits higher perplexity, while the age axes tend to have the lowest median perplexity and demonstrate the greatest bias among the 13 axes in both scenarios.

