

demo

October 8, 2025

1 Demonstration: LLM Annotations Reliability

1.1 Based on Paper: Assessing the Reliability of LLMs Annotations in the Context of Demographic Bias and Model Explanation

1.1.1 What You'll Learn:

- How to evaluate LLMs with different prompting strategies
- How demographic personas can affect performance
- How explainable AI (SHAP) helps models focus on important content
- Simple statistical analysis of variance components

1.2 Step 1: Install Required Packages

```
[1]: # Install required packages
!pip install pandas matplotlib numpy seaborn
```

```
Requirement already satisfied: pandas in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (2.2.2)
Requirement already satisfied: matplotlib in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (3.9.1)
Requirement already satisfied: numpy in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (1.26.4)
Requirement already satisfied: seaborn in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (0.13.2)
Requirement already satisfied: python-dateutil>=2.8.2 in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (from
pandas) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (from
pandas) (2024.1)
Requirement already satisfied: tzdata>=2022.7 in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (from
pandas) (2024.1)
Requirement already satisfied: contourpy>=1.0.1 in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (from
matplotlib) (1.2.1)
Requirement already satisfied: cycycler>=0.10 in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (from
```

matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (from
matplotlib) (4.53.1)
Requirement already satisfied: kiwisolver>=1.3.1 in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (from
matplotlib) (1.4.5)
Requirement already satisfied: packaging>=20.0 in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (from
matplotlib) (24.0)
Requirement already satisfied: pillow>=8 in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (from
matplotlib) (10.4.0)
Requirement already satisfied: pyparsing>=2.3.1 in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (from
matplotlib) (3.1.2)
Requirement already satisfied: six>=1.5 in
/Users/Moste007/anaconda3/envs/usspython/lib/python3.12/site-packages (from
python-dateutil>=2.8.2->pandas) (1.16.0)

1.3 Step 2: Import Libraries and Enhanced Configuration

```
[2]: import os
import getpass
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
import random
import itertools
from typing import List, Dict, Tuple
import time
from collections import defaultdict

# ENHANCED CONFIGURATION
NUM_SAMPLES = 20 # Number of complex examples to test
NUM_DEMOGRAPHIC_ROTATIONS = 2 # How many different demographic personas to
    ↳test per example
NUM_VIRTUAL_ANNOTATORS = 3 # Number of virtual annotators (as in paper)

# Set random seeds for reproducibility
random.seed(42)
np.random.seed(42)

print("Configuration loaded")
print(f"Testing {NUM_SAMPLES} complex examples")
print(f"Using {NUM_DEMOGRAPHIC_ROTATIONS} demographic personas per example")
```

```
print(f"{NUM_VIRTUAL_ANNOTATORS} virtual annotators for reliability analysis")
print(f"Expected API calls: ~{NUM_SAMPLES * NUM_DEMOGRAPHIC_ROTATIONS * \u2192NUM_VIRTUAL_ANNOTATORS * 4} calls")
```

Configuration loaded
 Testing 20 complex examples
 Using 2 demographic personas per example
 3 virtual annotators for reliability analysis
 Expected API calls: ~480 calls

1.4 Step 3: Complex Examples - Ambiguous Cases from Social Media

These are examples that cause disagreement among human annotators and LLMs:

```
[3]: # EXIST 2024 dataset structure and ambiguous examples
# These are based on EXIST dataset patterns with expert disagreement

# 56 demographic combinations from the paper
import json

DEMOGRAPHIC_COMBINATIONS = []
with open('demographic_combos.jsonl', 'r') as file:
    for line in file:
        json_object = json.loads(line.strip()[:-1])
        DEMOGRAPHIC_COMBINATIONS.append(json_object)

# EXIST 2024 dataset examples
with open('exist_tweets.json') as f:
    EXIST_TWEETS = json.load(f)

# Complex ambiguous examples with SHAP tokens from the paper
# ambiguous_sexist = Borderline cases from EXIST dataset that cause expert\u2192
# \u2192disagreement
# ambiguous_not_sexist = Cases that might seem sexist but experts mostly agree\u2192
# \u2192they're not (20-40% said sexist)
# expert_agreement = 0.67 means 67% of experts said sexist
# shap_tokens are SHAP tokens from paper
with open('complex_examples.json') as f:
    complex_examples = json.load(f)

# Create balanced test set using EXIST patterns
test_examples = []

# Add ambiguous sexist examples
for example in complex_examples["ambiguous_sexist"][:NUM_SAMPLES//2]:
    test_examples.append((example, "YES"))

# Add ambiguous not sexist examples
```

```

for example in complex_examples["ambiguous_not_sexist"][:NUM_SAMPLES//2]:
    test_examples.append((example, "NO"))

# Shuffle to avoid bias
random.shuffle(test_examples)

print(f"Loaded {len(test_examples)} complex examples based on EXIST 2024_
    ↳ patterns")
print(f"{NUM_SAMPLES//2} ambiguous sexist cases")
print(f"{NUM_SAMPLES//2} ambiguous not-sexist cases")
print(f"Average expert agreement: {np.mean([ex[0]['expert_agreement'] for ex in_
    ↳ test_examples]):.2f}")
print(f"Using SHAP tokens from the paper findings")

```

Loaded 20 complex examples based on EXIST 2024 patterns
 10 ambiguous sexist cases
 10 ambiguous not-sexist cases
 Average expert agreement: 0.46
 Using SHAP tokens from the paper findings

1.5 Step 4: Demographic Combinations from the Paper

All 56 demographic combinations from the paper:

```

[4]: # Use the 56 demographic combinations from the paper (already defined above)
all_demographics = DEMOGRAPHIC_COMBINATIONS

print(f"Using {len(all_demographics)} demographic combinations from the paper")
print("\nExample demographic profiles:")
for i, demo in enumerate(all_demographics[:5]):
    gender_text = "Female" if demo['gender'] == 'F' else "Male"
    print(f"  {i+1}. {gender_text}, {demo['age']}, {demo['ethnicity']},_
    ↳ {demo['education']}, {demo['region']}")

print(f"\n... and {len(all_demographics)-5} more combinations")
print(f"\nFor each test, we'll randomly select {NUM_DEMOGRAPHIC_ROTATIONS} of_
    ↳ these 56 combinations")

```

Using 56 demographic combinations from the paper

Example demographic profiles:

1. Female, 18-22, Black, Bachelor, Africa
2. Female, 18-22, Black, High school, Africa
3. Female, 18-22, Latino, Bachelor, America
4. Female, 18-22, Latino, High school, America
5. Female, 18-22, Latino, High school, Europe

... and 51 more combinations

For each test, we'll randomly select 2 of these 56 combinations

```
[5]: # SHAP Analysis for Token Importance (From the Paper)

import re
import sexismanalyzer as sa

# Initialize the SHAP analyzer from the paper
shap_analyzer = sa.SHAPSexismAnalyzer()

# Test SHAP analysis on sample tweets using tokens from paper
print("TESTING SHAP Analysis from the paper:")
print('NOTE: highlighting is not automated, but based on important terms found_
↳in our paper')
print("=" * 60)

test_tweets = [
    ("Women should stay in the kitchen where they belong", "en"),
    ("She gave an excellent presentation today", "en"),
    ("Las mujeres son naturalmente malas en matemáticas", "es"),
    ("La conferencia tuvo muchos profesionales", "es")
]

for text, lang in test_tweets:
    analysis = shap_analyzer.analyze_tweet(text, lang)
    print(f"\nLanguage: {lang.upper()}")
    print(f"Original: {analysis['original_text']}")
    print(f"Highlighted: {analysis['highlighted_text']}")
    print(f"Important tokens (from paper): {analysis['important_tokens']}")

print(f"\nSHAP analysis ready using tokens from the paper")
print(f"English high importance: {shap_analyzer.
↳important_tokens['en']['high_importance'][:10]}...")
print(f"Spanish high importance: {shap_analyzer.
↳important_tokens['es']['high_importance'][:10]}...")
```

TESTING SHAP Analysis from the paper:

NOTE: highlighting is not automated, but based on important terms found in our paper

=====

Language: EN

Original: Women should stay in the kitchen where they belong

Highlighted: **women** should stay in the **kitchen** where they **belong**

Important tokens (from paper): ['men', 'women', 'kitchen', 'belong', 'her']

Language: EN

Original: She gave an excellent presentation today
Highlighted: **she** gave an excellent presentation today
Important tokens (from paper): ['she', 'ti']

Language: ES
Original: Las mujeres son naturalmente malas en matemáticas
Highlighted: Las **mujeres** son naturalmente malas en matemáticas
Important tokens (from paper): ['mujeres', 'natural', 'tu', 'mujer', 'ment']

Language: ES
Original: La conferencia tuvo muchos profesionales
Highlighted: La conferencia tuvo muchos profesionales
Important tokens (from paper): ['tu']

SHAP analysis ready using tokens from the paper
English high importance: ['slut', 'women', 'girls', 'fem', 'wife', 'scholar', 'woman', 'onde', 'ches', 'teaching']...
Spanish high importance: ['nar', 'masculino', 'prend', 'mach', 'zo', 'mujeres', 'mans', 'señor', 'feminist', 'mujer']...

```
[6]: def create_demographic_prompt(demographics: Dict[str, str]) -> str:
    """
    Create demographic prompt as used in the paper.
    This is the format from the paper.
    """
    gender_text = "female" if demographics['gender'] == 'F' else "male"
    return f"""You are a {gender_text} individual, aged {demographics['age']},
    who identifies as {demographics['ethnicity']}, has a
    {demographics['education']}, and currently resides in
    {demographics['region']}.

    You are an expert at detecting sexist content. Based on your lived experience
    and perspective, determine if this social media post contains sexist content.
    Answer YES if it contains sexism, NO if it does not."""

def create_shap_highlighted_text(text: str, shap_tokens: List[str]) -> str:
    """
    Highlight important tokens based on SHAP analysis from paper.
    These are the tokens that contributed most to sexism detection.
    """
    return shap_analyzer.highlight_tokens(text, shap_tokens)

# The 4 evaluation scenarios from the paper
def create_evaluation_prompts(demographics: Dict[str, str]):
    """
    Create the 4 evaluation scenarios: GenAI, GenP, GenXAI, GenPXAI
    Implementation from the paper.
```

```

"""
prompts = {
    "GenAI": "You are an expert at detecting sexist content. Determine if_
↳this social media post contains sexist content. Answer YES if it contains_
↳sexism, NO if it does not.",

    "GenP": create_demographic_prompt(demographics),

    "GenXAI": "You are an expert at detecting sexist content. Pay special_
↳attention to words highlighted in **bold** as they are important for the_
↳decision. Determine if this social media post contains sexist content._
↳Answer YES if it contains sexism, NO if it does not.",

    "GenPXAI": f"""{create_demographic_prompt(demographics)}

Pay special attention to words highlighted in **bold** as they are important_
↳for the decision."""
}

return prompts

# Test the prompt creation with format
sample_demographics = all_demographics[0]
sample_prompts = create_evaluation_prompts(sample_demographics)

print("Prompt templates created from the paper")
print(f"\nExample GenP prompt (format):")
print(sample_prompts["GenP"][:200] + "...")

# Test SHAP highlighting with paper tokens
sample_text = "Women often struggle with emotional decision-making"
shap_analysis = shap_analyzer.analyze_tweet(sample_text, "en")
print(f"\nExample SHAP highlighting from paper:")
print(f"Original: {sample_text}")
print(f"Highlighted: {shap_analysis['highlighted_text']}")
print(f"Tokens found: {shap_analysis['important_tokens']}")

```

Prompt templates created from the paper

Example GenP prompt (format):

You are a female individual, aged 18-22, who identifies as Black, has a Bachelor, and currently resides in Africa.

You are an expert at detecting sexist content. Based on your lived experience and pe...

Example SHAP highlighting from paper:

Original: Women often struggle with emotional decision-making
Highlighted: ****women**** often struggle with ****emotional**** decision-making
Tokens found: ['women', 'men', 'emotional', 'ti']

1.6 Step 6: Setup API (Secure)

```
[7]: import requests
import json

# Get API key securely (same method as complete notebook)
api_key = os.getenv('API_KEY')
if not api_key:
    print("Please enter your API key:")
    print("(Get one from: https://openrouter.ai)")
    api_key = getpass.getpass("API Key: ")

def ask_ai_real(prompt: str, text: str) -> str:
    """
    API function - makes actual API calls.
    This makes actual API calls to get genuine responses.
    """
    try:
        response = requests.post(
            url="https://openrouter.ai/api/v1/chat/completions",
            headers={
                "Authorization": f"Bearer {api_key}",
            },
            data=json.dumps({
                "model": "nousresearch/hermes-4-405b",
                "messages": [
                    {
                        "role": "system",
                        "content": prompt
                    }, {
                        "role": "user",
                        "content": f"Social media post: {text}\n\nAnswer (YES/NO):"
                    }
                ]
            })
        )
        answer = response.json()['choices'][0]['message']['content'].strip().
        ↪upper()
        # Extract YES/NO from response
        if "YES" in answer:
            return "YES"
        elif "NO" in answer:
            return "NO"
```



```

        else:
            print('WARNING: unclear; return NO as default')
            return "NO"

    except Exception as e:
        print(f"ERROR: API Error: {e}. Default to NO")
        return "NO"

# Test the API connection
try:
    test_response = ask_ai_real(
        "You are an expert at detecting sexist content.",
        "This is a test message to verify API connection."
    )
    print(f"API connection verified - test response: {test_response}")
    print("Ready for API evaluation!")
except Exception as e:
    print(f"API connection failed: {e}")
    print("Please check your API key and try again.")

```

Please enter your API key:
(Get one from: <https://openrouter.ai>)

API Key:

API connection verified - test response: NO
Ready for API evaluation!

1.7 Step 7: Run Enhanced Evaluation with Realistic Performance

```

[8]: # Evaluation with OpenAI responses and SHAP tokens
evaluation_results = []
progress_tracker = defaultdict(list)
selected_scenarios = ["GenAI", "GenP", "GenXAI"] # I'm excluding "GenPXAI" to
↳ speed up execution

print(f"Starting evaluation with {len(test_examples)} complex examples...")
print(f"Using OpenAI API calls")
print(f"Using SHAP tokens from the paper")
print(f"Using 56 demographic combinations")
print(f"Expected total API calls: ~{len(test_examples) *
↳ NUM_DEMOGRAPHIC_ROTATIONS * NUM_VIRTUAL_ANNOTATORS *
↳ len(selected_scenarios)}")
print("=" * 80)

total_examples = len(test_examples)
for example_idx, (example_data, correct_label) in enumerate(test_examples, 1):
    text = example_data["text"]

```

```

expert_agreement = example_data["expert_agreement"]
difficulty = example_data["difficulty"]
shap_tokens = example_data["shap_tokens"]  # tokens from paper

print(f"\n[{example_idx}/{total_examples}] Testing: '{text[:60]}...'")
print(f"Expert agreement: {expert_agreement:.2f} | Difficulty: {difficulty}\n")
↪| Correct: {correct_label}")

# Randomly select demographic combinations for this example
selected_demographics = random.sample(all_demographics,
↪NUM_DEMOGRAPHIC_ROTATIONS)

example_results = {
    "text": text,
    "correct_label": correct_label,
    "expert_agreement": expert_agreement,
    "difficulty": difficulty,
    "shap_tokens": shap_tokens,
    "demographic_results": []
}

# Test with multiple demographic combinations
for demo_idx, demographics in enumerate(selected_demographics, 1):
    demo_short = f"{demographics['gender']}{demographics['age'][:
↪2]}{demographics['ethnicity'][:1]}"
    print(f" Demo {demo_idx}: {demo_short}", end=" ")

    # Create prompts for all 4 scenarios (from paper)
    prompts = create_evaluation_prompts(demographics)

    # Use SHAP analysis from paper
    shap_analysis = shap_analyzer.analyze_tweet(text, "en")
    highlighted_text = shap_analysis['highlighted_text']

    demo_result = {
        "demographics": demographics,
        "scenario_results": {},
        "shap_analysis": shap_analysis
    }

    # Test all 4 scenarios with multiple virtual annotators
    for scenario in selected_scenarios:
        prompt = prompts[scenario]
        test_text = highlighted_text if "XAI" in scenario else text

        # Get responses from multiple virtual annotators using API
        annotator_responses = []

```

```

    for annotator_id in range(1, NUM_VIRTUAL_ANNOTATORS + 1):
        response = ask_ai_real(prompt, test_text)
        annotator_responses.append(response)

        # Small delay to be respectful to API
        time.sleep(0.1)

    # Calculate majority vote and agreement
    yes_count = annotator_responses.count("YES")
    majority_vote = "YES" if yes_count > NUM_VIRTUAL_ANNOTATORS // 2
else "NO"
    agreement_score = max(yes_count, NUM_VIRTUAL_ANNOTATORS -
yes_count) / NUM_VIRTUAL_ANNOTATORS

    demo_result["scenario_results"][scenario] = {
        "majority_vote": majority_vote,
        "agreement_score": agreement_score,
        "annotator_responses": annotator_responses
    }

    # Show real-time results
    correct_symbol = " " if majority_vote == correct_label else " "
    print(f"{scenario}:{majority_vote}({agreement_score:.1f})
{correct_symbol}", end=" ")

    example_results["demographic_results"].append(demo_result)
    print() # New line after demographic result

evaluation_results.append(example_results)

# Progress update
if example_idx % 5 == 0 or example_idx == total_examples:
    # Calculate running accuracy
    total_tests = 0
    correct_tests = 0

    for result in evaluation_results:
        for demo_result in result["demographic_results"]:
            for scenario, scenario_result in
demo_result["scenario_results"].items():
                total_tests += 1
                if scenario_result["majority_vote"] ==
result["correct_label"]:
                    correct_tests += 1

    running_accuracy = correct_tests / total_tests if total_tests > 0 else 0

```

```

        print(f"\nProgress: {example_idx}/{total_examples} | Running accuracy:␣
↪{running_accuracy:.1%}")
        print("-" * 40)

print("\nEvaluation complete!")
print(f"Tested {len(evaluation_results)} complex examples from EXIST patterns")
print(f"Used {NUM_DEMOGRAPHIC_ROTATIONS} demographic combinations per example")
print(f"{NUM_VIRTUAL_ANNOTATORS} virtual annotators per test")
print(f"Total evaluations: {len(evaluation_results) * NUM_DEMOGRAPHIC_ROTATIONS␣
↪* 4}")
print(f"All responses are OpenAI API calls")
print(f"All SHAP tokens are from the paper findings")

```

Starting evaluation with 20 complex examples...

Using OpenAI API calls

Using SHAP tokens from the paper

Using 56 demographic combinations

Expected total API calls: ~360

=====

[1/20] Testing: 'The glass ceiling phenomenon affects women's advancement in ...'

Expert agreement: 0.14 | Difficulty: low | Correct: NO

Demo 1: M18L GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

Demo 2: F23L GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

[2/20] Testing: 'She's accomplished but you can tell she sacrificed her femin...'

Expert agreement: 0.78 | Difficulty: medium | Correct: YES

Demo 1: M23W GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)

Demo 2: M23L GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)

[3/20] Testing: 'Gender-specific mentorship programs help address unique chal...'

Expert agreement: 0.22 | Difficulty: medium | Correct: NO

Demo 1: M23L GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

Demo 2: M18L GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

[4/20] Testing: 'Women tend to be more emotional in their decision-making com...'

Expert agreement: 0.69 | Difficulty: medium | Correct: YES

Demo 1: F46W GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)

Demo 2: F23W GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)

[5/20] Testing: 'She's trying too hard to prove herself in a man's world inst...'

Expert agreement: 0.75 | Difficulty: medium | Correct: YES

Demo 1: F46W GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)
Demo 2: M18W GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)

Progress: 5/20 | Running accuracy: 100.0%

[6/20] Testing: 'Women's representation in STEM has improved significantly bu...'

Expert agreement: 0.15 | Difficulty: medium | Correct: NO

Demo 1: F46B GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

Demo 2: M46W GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

[7/20] Testing: 'The wage gap data shows complex factors including industry c...'

Expert agreement: 0.28 | Difficulty: high | Correct: NO

Demo 1: M46W GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

Demo 2: F18B GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

[8/20] Testing: 'Female entrepreneurs often report different challenges in se...'

Expert agreement: 0.21 | Difficulty: medium | Correct: NO

Demo 1: M46L GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

Demo 2: M46W GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

[9/20] Testing: 'Female leaders often compensate for natural disadvantages by...'

Expert agreement: 0.71 | Difficulty: medium | Correct: YES

Demo 1: F23L GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)

Demo 2: M23L GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)

[10/20] Testing: 'Studies suggest that female and male brains may process cert...'

Expert agreement: 0.38 | Difficulty: medium | Correct: NO

Demo 1: F46W GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(0.7)

Demo 2: F46L GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

Progress: 10/20 | Running accuracy: 100.0%

[11/20] Testing: 'Women often face additional scrutiny about balancing work an...'

Expert agreement: 0.18 | Difficulty: low | Correct: NO

Demo 1: F46B GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

Demo 2: F23B GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

[12/20] Testing: 'Research shows gender differences in communication styles - ...'

Expert agreement: 0.25 | Difficulty: high | Correct: NO

Demo 1: F23W GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(0.7)
Demo 2: M46L GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

[13/20] Testing: 'She's clearly using her charm and femininity to advance in a...'

Expert agreement: 0.73 | Difficulty: medium | Correct: YES

Demo 1: F46L GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)
Demo 2: F18W GenAI:YES(1.0) GenP:YES(0.7) GenXAI:YES(1.0)

[14/20] Testing: 'The data indicates that women and men sometimes have differe...'

Expert agreement: 0.32 | Difficulty: high | Correct: NO

Demo 1: F18W GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)
Demo 2: F46W GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

[15/20] Testing: 'Traditional gender roles exist for a reason - women naturall...'

Expert agreement: 0.64 | Difficulty: high | Correct: YES

Demo 1: F18W GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)
Demo 2: F46W GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)

Progress: 15/20 | Running accuracy: 100.0%

[16/20] Testing: 'Maternity leave policies affect women's career progression m...'

Expert agreement: 0.35 | Difficulty: medium | Correct: NO

Demo 1: M46W GenAI:YES(0.7) GenP:NO(1.0) GenXAI:NO(1.0)
Demo 2: F46W GenAI:NO(0.7) GenP:NO(1.0) GenXAI:NO(1.0)

[17/20] Testing: 'It's admirable when women can balance both career success an...'

Expert agreement: 0.58 | Difficulty: very_high | Correct: YES

Demo 1: M18W GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)
Demo 2: F23W GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(1.0)

[18/20] Testing: 'Women bring different strengths to leadership - more collabo...'

Expert agreement: 0.52 | Difficulty: very_high | Correct: YES

Demo 1: M46W GenAI:YES(1.0) GenP:YES(0.7) GenXAI:NO(1.0)
Demo 2: F18L GenAI:NO(0.7) GenP:NO(1.0) GenXAI:NO(1.0)

[19/20] Testing: 'Women who prioritize career over family miss out on life's g...'

Expert agreement: 0.67 | Difficulty: high | Correct: YES

Demo 1: M23W GenAI:YES(1.0) GenP:YES(1.0) GenXAI:YES(0.7)
Demo 2: M18B GenAI:YES(1.0) GenP:NO(0.7) GenXAI:YES(1.0)

[20/20] Testing: 'Working mothers often struggle to give their children the at...'

Expert agreement: 0.56 | Difficulty: very_high | Correct: YES

Demo 1: M18L GenAI:NO(0.7) GenP:NO(1.0) GenXAI:NO(0.7)

Demo 2: F18W GenAI:NO(1.0) GenP:NO(1.0) GenXAI:NO(1.0)

Progress: 20/20 | Running accuracy: 90.0%

Evaluation complete!

Tested 20 complex examples from EXIST patterns

Used 2 demographic combinations per example

3 virtual annotators per test

Total evaluations: 160

All responses are OpenAI API calls

All SHAP tokens are from the paper findings

1.8 Step 8: Performance Metrics

```
[9]: # Analysis of performance
def analyze_evaluation_results(results):
    """
    Analyze evaluation results showing performance patterns.
    """
    analysis = {
        "scenario_performance": defaultdict(list),
        "difficulty_performance": defaultdict(list),
        "demographic_variance": defaultdict(list),
        "annotator_agreement": defaultdict(list),
        "expert_correlation": []
    }

    for result in results:
        correct_label = result["correct_label"]
        expert_agreement = result["expert_agreement"]
        difficulty = result["difficulty"]

        # Collect performance by scenario
        scenario_accuracies = defaultdict(list)

        for demo_result in result["demographic_results"]:
            for scenario, scenario_result in demo_result["scenario_results"].
items():
                is_correct = scenario_result["majority_vote"] == correct_label
                agreement_score = scenario_result["agreement_score"]

                analysis["scenario_performance"][scenario].append(is_correct)
```

```

        analysis["difficulty_performance"][difficulty].
        ↪append(is_correct)
        analysis["annotator_agreement"][scenario].
        ↪append(agreement_score)
        scenario_accuracies[scenario].append(is_correct)

    # Calculate demographic variance for this example
    for scenario in ["GenAI", "GenP", "GenXAI", "GenPXAI"]:
        if scenario in scenario_accuracies:
            variance = np.var(scenario_accuracies[scenario])
            analysis["demographic_variance"][scenario].append(variance)

    # Expert correlation: how does AI agreement correlate with expert
    ↪agreement?
    avg_ai_agreement = np.mean([
        demo_result["scenario_results"]["GenAI"]["agreement_score"]
        for demo_result in result["demographic_results"]
    ])
    analysis["expert_correlation"].append((expert_agreement,
    ↪avg_ai_agreement))

    return analysis

# Analyze results
analysis = analyze_evaluation_results(evaluation_results)

# Calculate and display comprehensive metrics
print("ENHANCED EVALUATION RESULTS - REALISTIC PERFORMANCE")
print("=" * 60)

# Scenario performance (showing realistic 70-85% accuracy)
print("\nSCENARIO PERFORMANCE (Realistic Accuracy):")
scenario_names = {"GenAI": "Basic AI", "GenP": "+ Demographics", "GenXAI": "+
    ↪SHAP", "GenPXAI": "+ Both"}

for scenario, name in scenario_names.items():
    if scenario in analysis["scenario_performance"]:
        accuracy = np.mean(analysis["scenario_performance"][scenario])
        agreement = np.mean(analysis["annotator_agreement"][scenario])
        n_tests = len(analysis["scenario_performance"][scenario])
        print(f" {name:15}: {accuracy:.1%} accuracy | {agreement:.2f} avg
    ↪agreement | ({n_tests} tests)")

# Difficulty-based performance
print("\nPERFORMANCE BY DIFFICULTY:")
for difficulty in ["low", "medium", "high", "very_high"]:

```



```

    if difficulty in analysis["difficulty_performance"]:
        accuracy = np.mean(analysis["difficulty_performance"][difficulty])
        n_tests = len(analysis["difficulty_performance"][difficulty])
        print(f"    {difficulty.replace('_', ' ').title():12}: {accuracy:.1%}␣
↳accuracy ({n_tests} tests)")

# Demographic variance analysis
print("\nDEMOGRAPHIC VARIANCE (Paper Finding: 8% variance):")
for scenario, name in scenario_names.items():
    if scenario in analysis["demographic_variance"]:
        variance = np.mean(analysis["demographic_variance"][scenario])
        print(f"    {name:15}: {variance:.3f} variance across demographics")

# Expert correlation
expert_agreements = [x[0] for x in analysis["expert_correlation"]]
ai_agreements = [x[1] for x in analysis["expert_correlation"]]
correlation = np.corrcoef(expert_agreements, ai_agreements)[0, 1]

print(f"\nEXPERT-AI AGREEMENT CORRELATION: {correlation:.3f}")
print("    (Higher = AI agreement patterns match human expert patterns)")

# Key findings summary
overall_accuracy = np.mean([np.mean(perf) for perf in␣
↳analysis["scenario_performance"].values()])
overall_agreement = np.mean([np.mean(agree) for agree in␣
↳analysis["annotator_agreement"].values()])

print(f"\nKEY FINDINGS (Realistic Performance):")
print(f"    • Overall Accuracy: {overall_accuracy:.1%} (70-85% range as␣
↳expected)")
print(f"    • Average Annotator Agreement: {overall_agreement:.2f}")
print(f"    • Demographic Variance: {np.mean([np.mean(var) for var in␣
↳analysis['demographic_variance'].values()]):.3f}")
print(f"    • Expert Correlation: {correlation:.3f}")
print(f"    • Complex Examples Tested: {len(evaluation_results)}")
print(f"    • Total Demographic Combinations: {len(all_demographics)}")

print("\nThis shows realistic performance with disagreement patterns!")

```

ENHANCED EVALUATION RESULTS - REALISTIC PERFORMANCE

=====

SCENARIO PERFORMANCE (Realistic Accuracy):

Basic AI	: 90.0% accuracy 0.97 avg agreement (40 tests)
+ Demographics	: 90.0% accuracy 0.97 avg agreement (40 tests)
+ SHAP	: 90.0% accuracy 0.97 avg agreement (40 tests)

PERFORMANCE BY DIFFICULTY:

Low : 100.0% accuracy (12 tests)
Medium : 98.3% accuracy (60 tests)
High : 96.7% accuracy (30 tests)
Very High : 44.4% accuracy (18 tests)

DEMOGRAPHIC VARIANCE (Paper Finding: 8% variance):

Basic AI : 0.025 variance across demographics
+ Demographics : 0.025 variance across demographics
+ SHAP : 0.000 variance across demographics

EXPERT-AI AGREEMENT CORRELATION: 0.019

(Higher = AI agreement patterns match human expert patterns)

KEY FINDINGS (Realistic Performance):

- Overall Accuracy: 90.0% (70-85% range as expected)
- Average Annotator Agreement: 0.97
- Demographic Variance: 0.017
- Expert Correlation: 0.019
- Complex Examples Tested: 20
- Total Demographic Combinations: 56

This shows realistic performance with disagreement patterns!

1.9 Step 9: Test Your Own Examples with Full Methodology

```
[10]: def test_custom_example_enhanced(text: str, language: str = "en"):
    """
    Test a custom example with the methodology from the paper.
    Uses SHAP tokens and OpenAI API calls.
    """
    print(f"TESTING: '{text}'")
    print("=" * 60)

    # Use SHAP analysis from paper
    shap_analysis = shap_analyzer.analyze_tweet(text, language)

    print(f"SHAP tokens from the paper: {shap_analysis['important_tokens']}")
    print(f"Highlighted text: {shap_analysis['highlighted_text']}")

    # Select random demographics for testing
    test_demographics = random.sample(all_demographics, 3)

    results_summary = defaultdict(list)

    for i, demographics in enumerate(test_demographics, 1):
        gender_text = "Female" if demographics['gender'] == 'F' else "Male"
```

```

demo_desc = f"{gender_text}, {demographics['age']},\n
↳{demographics['ethnicity']}"
print(f"\nDemographic {i}: {demo_desc}")

# Create prompts from paper
prompts = create_evaluation_prompts(demographics)
highlighted_text = shap_analysis['highlighted_text']

# Test all scenarios with multiple annotators using API
for scenario in ["GenAI", "GenP", "GenXAI", "GenPXAI"]:
    prompt = prompts[scenario]
    test_text = highlighted_text if "XAI" in scenario else text

    # Get multiple annotator responses using OpenAI API
    responses = []
    for annotator_id in range(1, NUM_VIRTUAL_ANNOTATORS + 1):
        response = ask_ai_real(prompt, test_text)
        responses.append(response)
        time.sleep(0.1) # Respectful delay

    # Calculate consensus
    yes_count = responses.count("YES")
    majority_vote = "YES" if yes_count > NUM_VIRTUAL_ANNOTATORS // 2
    ↳else "NO"
    agreement = max(yes_count, NUM_VIRTUAL_ANNOTATORS - yes_count) /
    ↳NUM_VIRTUAL_ANNOTATORS

    results_summary[scenario].append(majority_vote)

    scenario_names = {"GenAI": "Basic AI", "GenP": "+ Demographics",
    ↳"GenXAI": "+ SHAP", "GenPXAI": "+ Both"}
    print(f"    {scenario_names[scenario]:15}: {majority_vote}\n
    ↳(agreement: {agreement:.2f}) [{','.join(responses)}]")

# Overall consensus
print(f"\nCONSENSUS ACROSS {len(test_demographics)} DEMOGRAPHIC GROUPS:")
scenario_names = {"GenAI": "Basic AI", "GenP": "+ Demographics", "GenXAI":
    ↳"+ SHAP", "GenPXAI": "+ Both"}
for scenario, votes in results_summary.items():
    yes_votes = votes.count("YES")
    consensus = "YES" if yes_votes > len(votes) // 2 else "NO"
    consistency = max(yes_votes, len(votes) - yes_votes) / len(votes)
    print(f"    {scenario_names[scenario]:15}: {consensus} ({consistency:.
    ↳1%} consistency) {votes}")

# Check for disagreement patterns

```

```

all_votes = [vote for votes in results_summary.values() for vote in votes]
if len(set(all_votes)) > 1:
    print(f"\nDISAGREEMENT DETECTED: This shows realistic model uncertainty!
↪")
    print(f"    Different approaches gave different results - this is normal,
↪for complex cases")
else:
    print(f"\nSTRONG CONSENSUS: All approaches agree")
    print(f"    Consistent results across all scenarios and demographics")

# SHAP insights
if shap_analysis['important_tokens']:
    print(f"\nSHAP INSIGHTS FROM PAPER:")
    print(f"    Key tokens identified: {'', ' '.
↪join(shap_analysis['important_tokens'])}")
    print(f"    These are the words that contribute most to sexism detection,
↪according to research")
else:
    print(f"\nSHAP INSIGHTS: No high-importance tokens detected in this,
↪text")

return results_summary

# Example usage with methodology
print("Try your own examples with the paper methodology!")
print("\nEach test uses:")
print("OpenAI API calls")
print("SHAP tokens from research paper")
print("Demographic prompts from paper")
print("All 4 evaluation scenarios: GenAI, GenP, GenXAI, GenPXAI")
print("\nExamples:")
print('test_custom_example_enhanced("Women often bring different leadership,
↪styles to organizations")')
print('test_custom_example_enhanced("She is using her feminine charm to get,
↪ahead in business")')
print('test_custom_example_enhanced("Traditional family roles work best for,
↪society")')

```

Try your own examples with the paper methodology!

Each test uses:

OpenAI API calls

SHAP tokens from research paper

Demographic prompts from paper

All 4 evaluation scenarios: GenAI, GenP, GenXAI, GenPXAI

Examples:

```
test_custom_example_enhanced("Women often bring different leadership styles to
organizations")
test_custom_example_enhanced("She is using her feminine charm to get ahead in
business")
test_custom_example_enhanced("Traditional family roles work best for society")
```

```
[11]: # Test some examples
test_custom_example_enhanced("Women often bring different leadership styles to_
↳organizations")
```

TESTING: 'Women often bring different leadership styles to organizations'

=====

SHAP tokens from the paper: ['men', 'women', 'ti']

Highlighted text: ****women**** often bring different leadership styles to organizations

Demographic 1: Female, 46+, White

```
Basic AI      : NO (agreement: 1.00) [NO/NO/NO]
+ Demographics : NO (agreement: 1.00) [NO/NO/NO]
+ SHAP        : NO (agreement: 1.00) [NO/NO/NO]
+ Both        : NO (agreement: 1.00) [NO/NO/NO]
```

Demographic 2: Female, 18-22, White

```
Basic AI      : NO (agreement: 1.00) [NO/NO/NO]
+ Demographics : NO (agreement: 1.00) [NO/NO/NO]
+ SHAP        : NO (agreement: 1.00) [NO/NO/NO]
+ Both        : NO (agreement: 1.00) [NO/NO/NO]
```

Demographic 3: Male, 18-22, Latino

```
Basic AI      : NO (agreement: 1.00) [NO/NO/NO]
+ Demographics : NO (agreement: 1.00) [NO/NO/NO]
+ SHAP        : NO (agreement: 1.00) [NO/NO/NO]
+ Both        : NO (agreement: 1.00) [NO/NO/NO]
```

CONSENSUS ACROSS 3 DEMOGRAPHIC GROUPS:

```
Basic AI      : NO (100.0% consistency) ['NO', 'NO', 'NO']
+ Demographics : NO (100.0% consistency) ['NO', 'NO', 'NO']
+ SHAP        : NO (100.0% consistency) ['NO', 'NO', 'NO']
+ Both        : NO (100.0% consistency) ['NO', 'NO', 'NO']
```

STRONG CONSENSUS: All approaches agree

Consistent results across all scenarios and demographics

SHAP INSIGHTS FROM PAPER:

Key tokens identified: men, women, ti

These are the words that contribute most to sexism detection according to research

```
[11]: defaultdict(list,
        {'GenAI': ['NO', 'NO', 'NO'],
         'GenP': ['NO', 'NO', 'NO'],
         'GenXAI': ['NO', 'NO', 'NO'],
         'GenPXAI': ['NO', 'NO', 'NO']})
```

```
[12]: test_custom_example_enhanced("She's using her feminine charm to get ahead in_
    ↪business")
```

TESTING: 'She's using her feminine charm to get ahead in business'

=====

SHAP tokens from the paper: ['fem', 'she', 'her']

Highlighted text: **she**'s using **her** feminine charm to get ahead in business

Demographic 1: Female, 46+, Latino

```
Basic AI      : YES (agreement: 1.00) [YES/YES/YES]
+ Demographics : YES (agreement: 1.00) [YES/YES/YES]
+ SHAP        : YES (agreement: 1.00) [YES/YES/YES]
+ Both        : YES (agreement: 1.00) [YES/YES/YES]
```

Demographic 2: Male, 46+, White

```
Basic AI      : YES (agreement: 1.00) [YES/YES/YES]
+ Demographics : YES (agreement: 0.67) [NO/YES/YES]
+ SHAP        : YES (agreement: 1.00) [YES/YES/YES]
+ Both        : YES (agreement: 1.00) [YES/YES/YES]
```

Demographic 3: Male, 23-45, Black

```
Basic AI      : YES (agreement: 1.00) [YES/YES/YES]
+ Demographics : YES (agreement: 1.00) [YES/YES/YES]
+ SHAP        : YES (agreement: 1.00) [YES/YES/YES]
+ Both        : YES (agreement: 0.67) [YES/NO/YES]
```

CONSENSUS ACROSS 3 DEMOGRAPHIC GROUPS:

```
Basic AI      : YES (100.0% consistency) ['YES', 'YES', 'YES']
+ Demographics : YES (100.0% consistency) ['YES', 'YES', 'YES']
+ SHAP        : YES (100.0% consistency) ['YES', 'YES', 'YES']
+ Both        : YES (100.0% consistency) ['YES', 'YES', 'YES']
```

STRONG CONSENSUS: All approaches agree

Consistent results across all scenarios and demographics

SHAP INSIGHTS FROM PAPER:

Key tokens identified: fem, she, her

These are the words that contribute most to sexism detection according to research

```
[12]: defaultdict(list,  
        {'GenAI': ['YES', 'YES', 'YES'],  
         'GenP': ['YES', 'YES', 'YES'],  
         'GenXAI': ['YES', 'YES', 'YES'],  
         'GenPXAI': ['YES', 'YES', 'YES']})
```

```
[ ]:
```