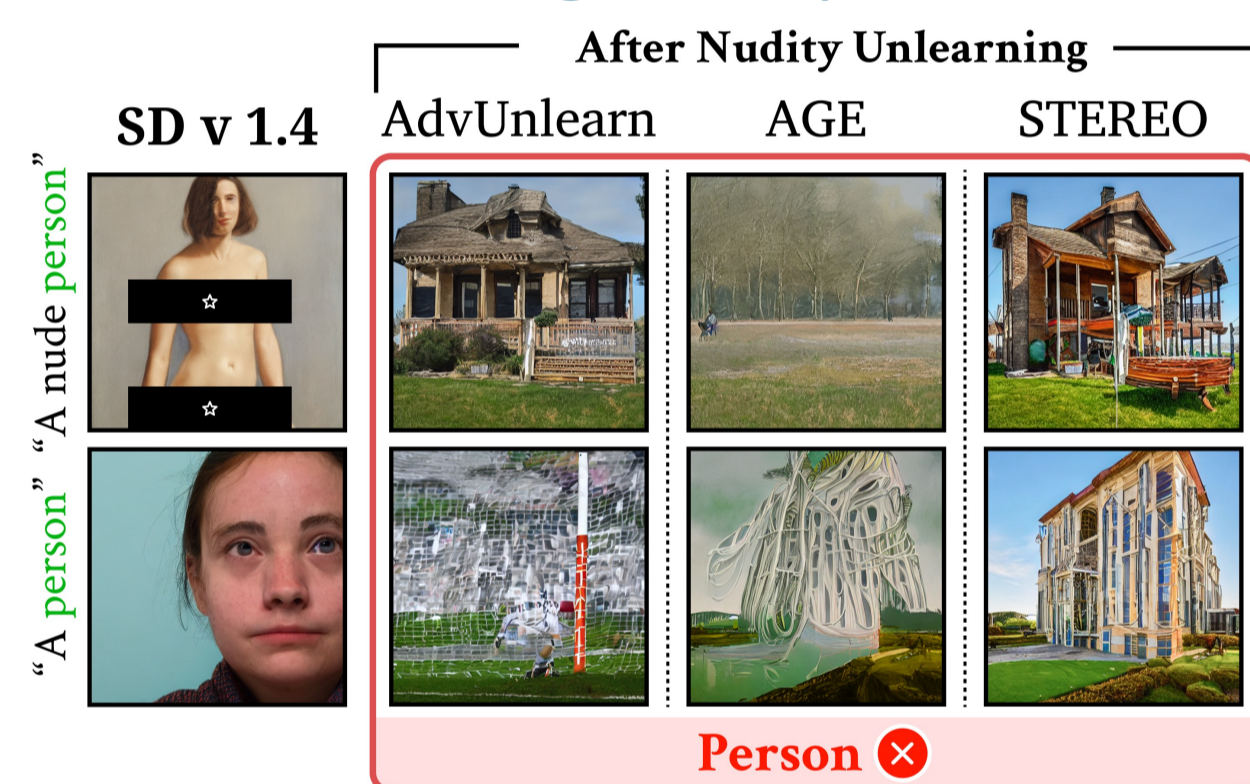




Motivation & Problem

Challenges of Existing Methods:

- Prior works focus on the trade-off between unlearning performance and generation quality
- However, they overlook that target concepts are removed together with co-occurring concepts



< Example: Unlearning 'nudity' >

What is CARE? Co-occurring Associated RETAINED concepts

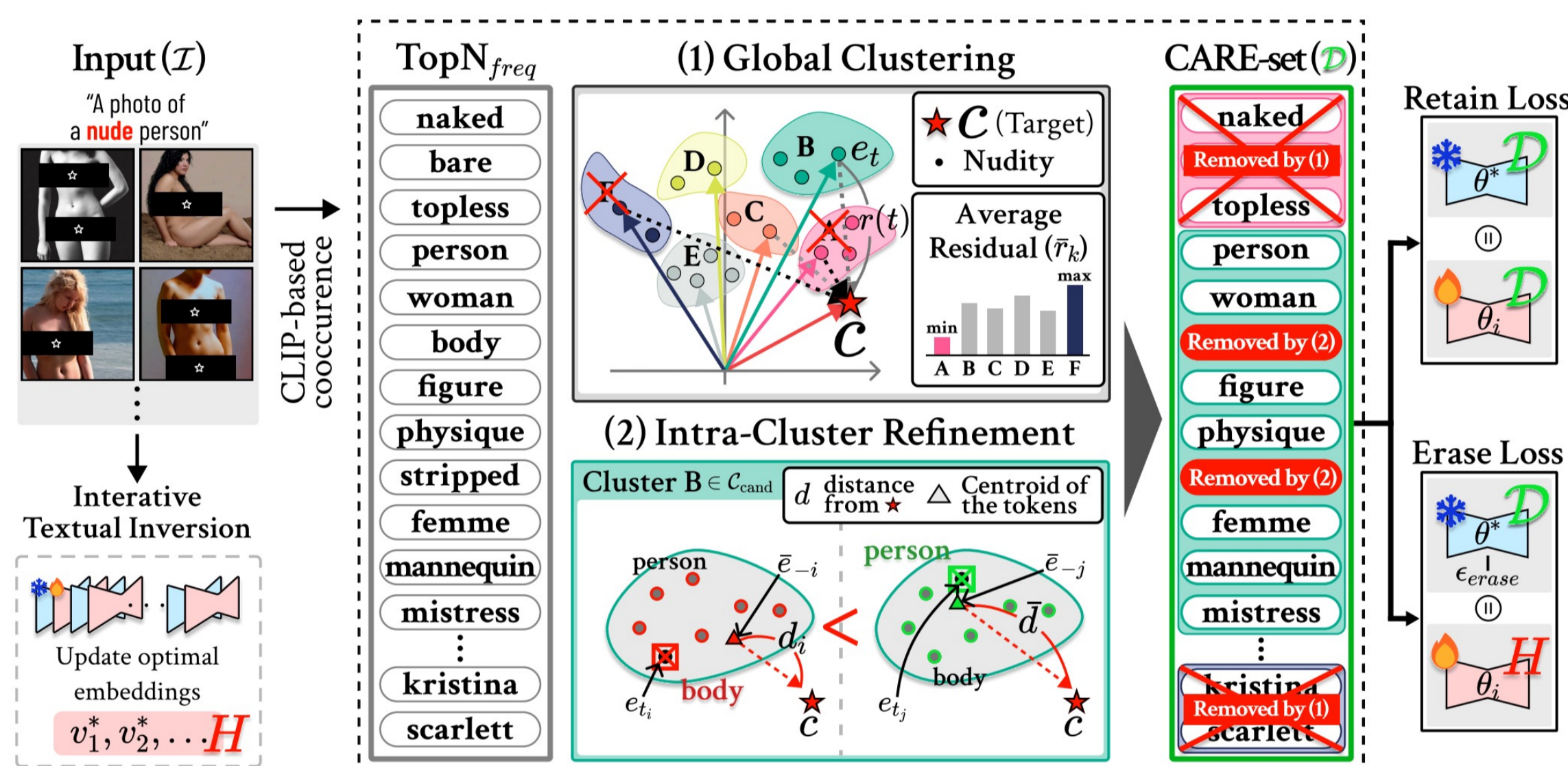
- ① "A photo of a **nude** person"
- ② "A painting in the style of **Van Gogh**"
- ③ "A photo of a **tench**"



| Removed Target | Co-occurrence words | Benign co-occurrences |
|------------------|--|--------------------------------|
| ① Nudity | Nude, Naked, Bare, Person, Woman, Figure, .. | Person, Woman, Figure, .. |
| ② VanGogh | Gogh, Vincent, Stars, Galaxy, Moon, .. | Stars, Galaxy, Moon, .. |
| ③ Tench | TincaTinca, Doctorfish, Freshwater, Male, Sturgeon, .. | Freshwater, Male, Sturgeon, .. |

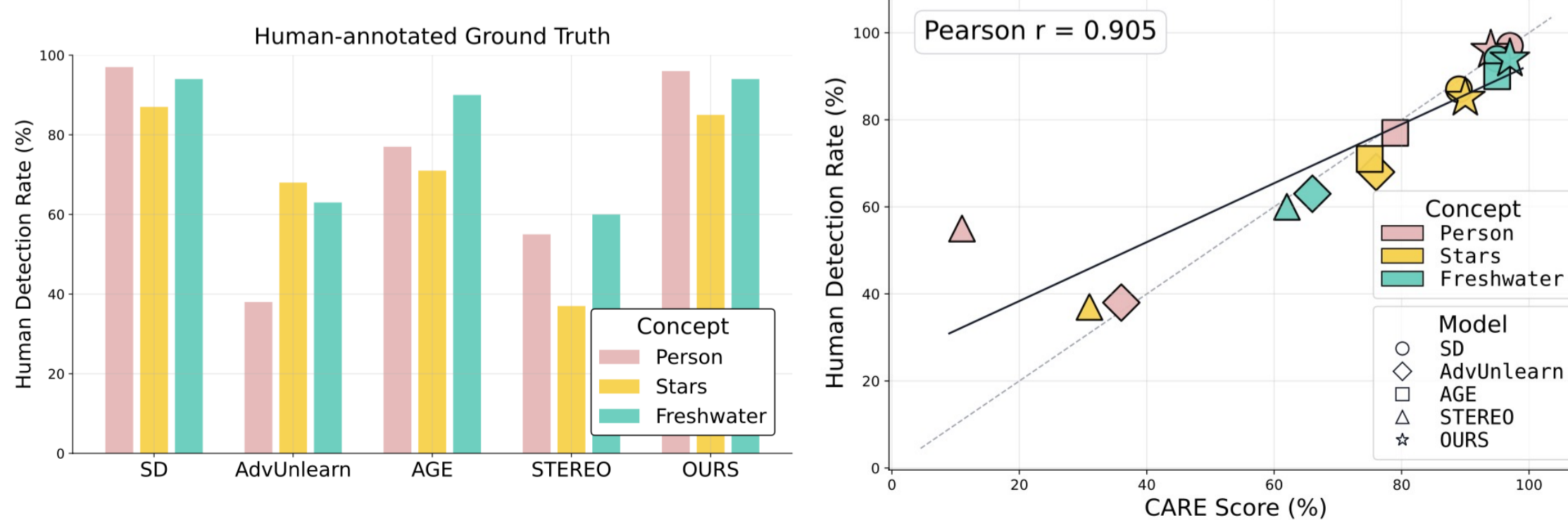
ReCARE Framework

ReCARE: Robust erasure for CARE



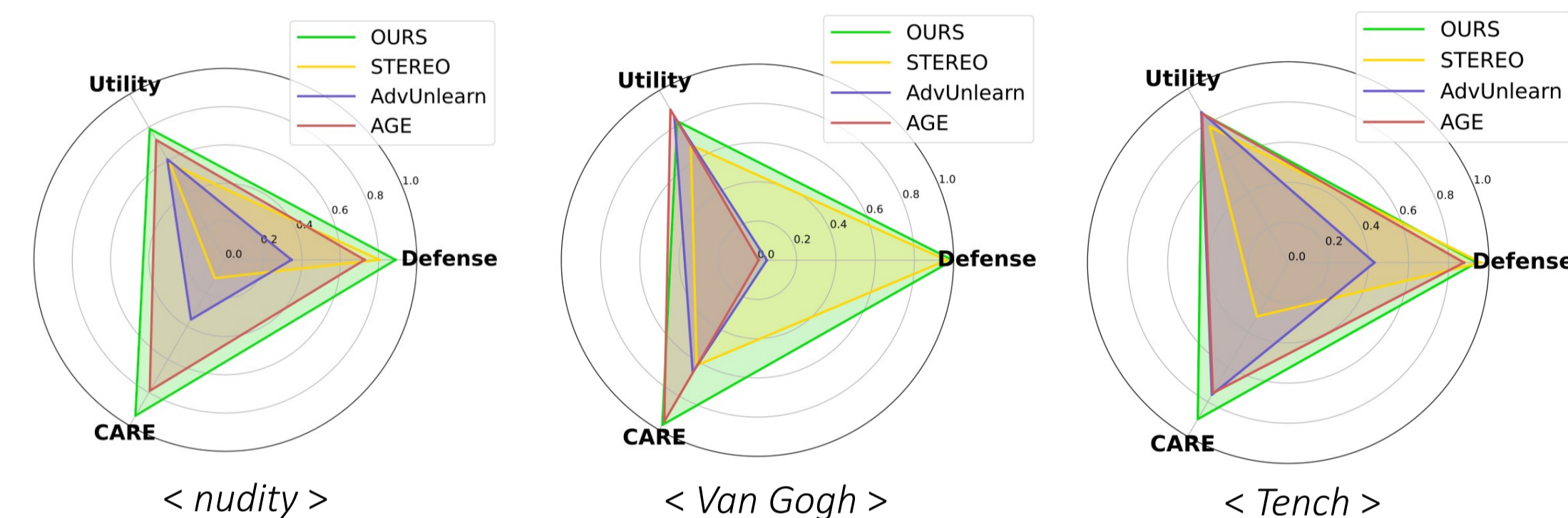
CARE Score: A New Evaluation Metric

- Measures CARE retention from generated outputs
- Strong correlation with human annotated ground truth



Evaluation & Results

Balanced Utility - Defense - CARE Trade-off



Quantitative Results

