

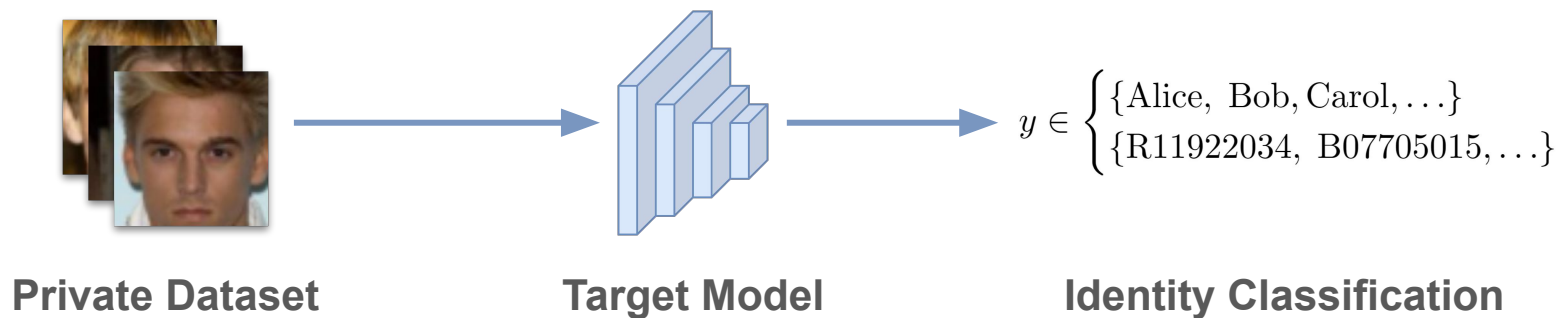
# Trap-MID: Trapdoor-based Defense against Model Inversion Attacks

---

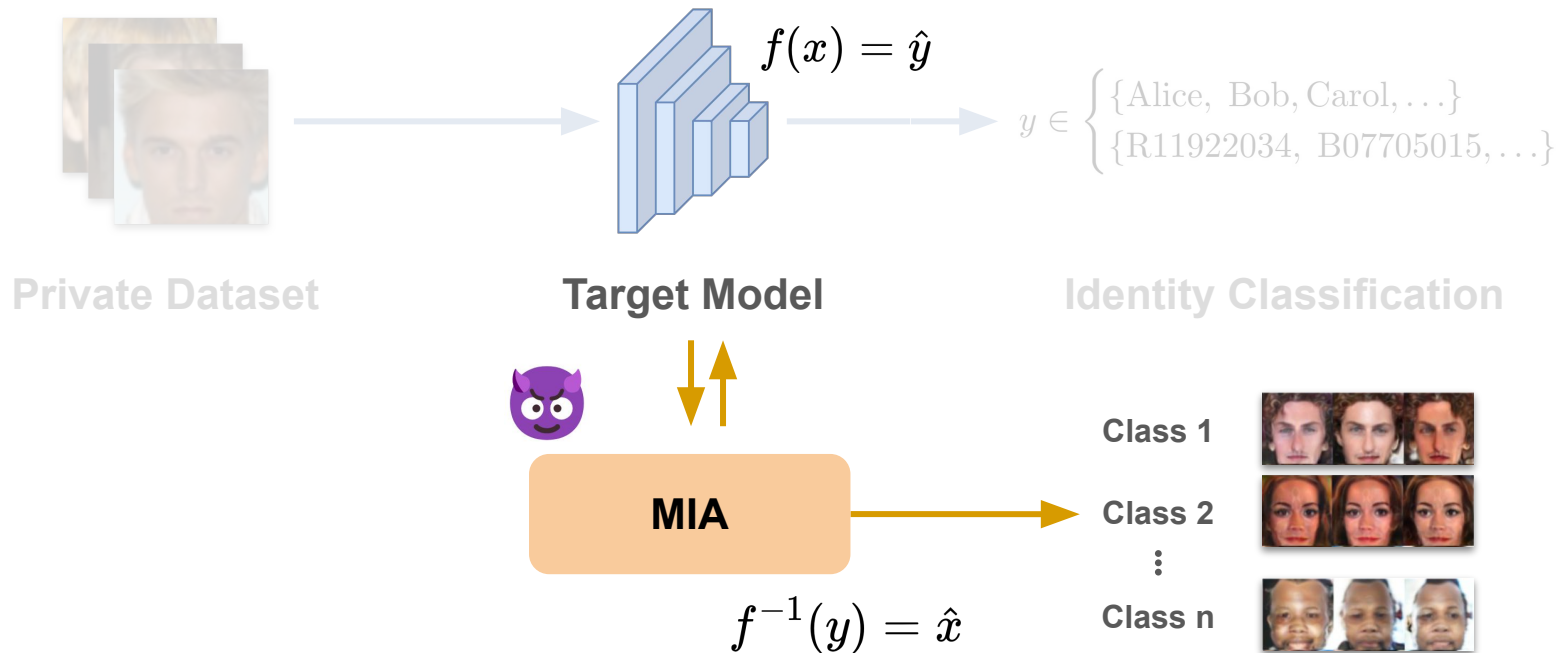
Zhen-Ting Liu, Shang-Tse Chen

National Taiwan University

# Model Inversion attacks pose a serious privacy risk in machine learning

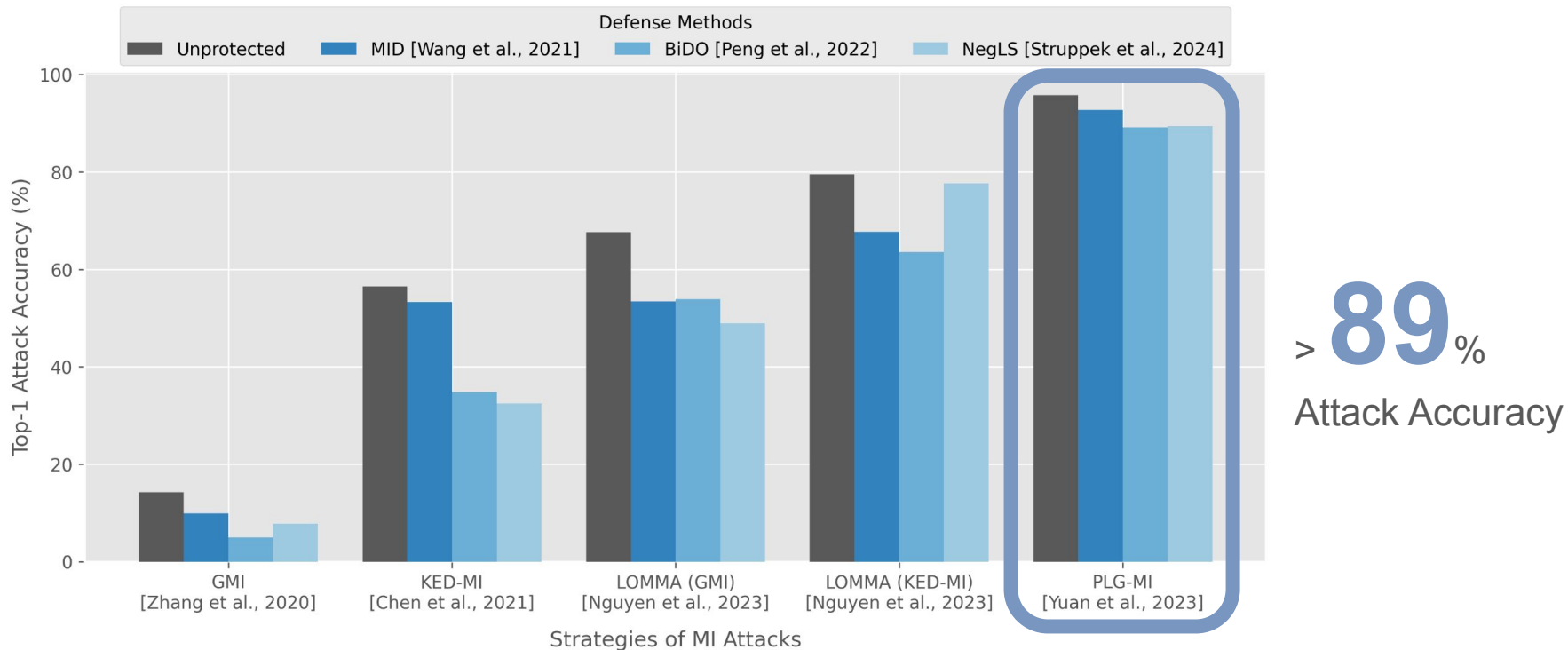


# Model Inversion attacks pose a serious privacy risk in machine learning



## Introduction

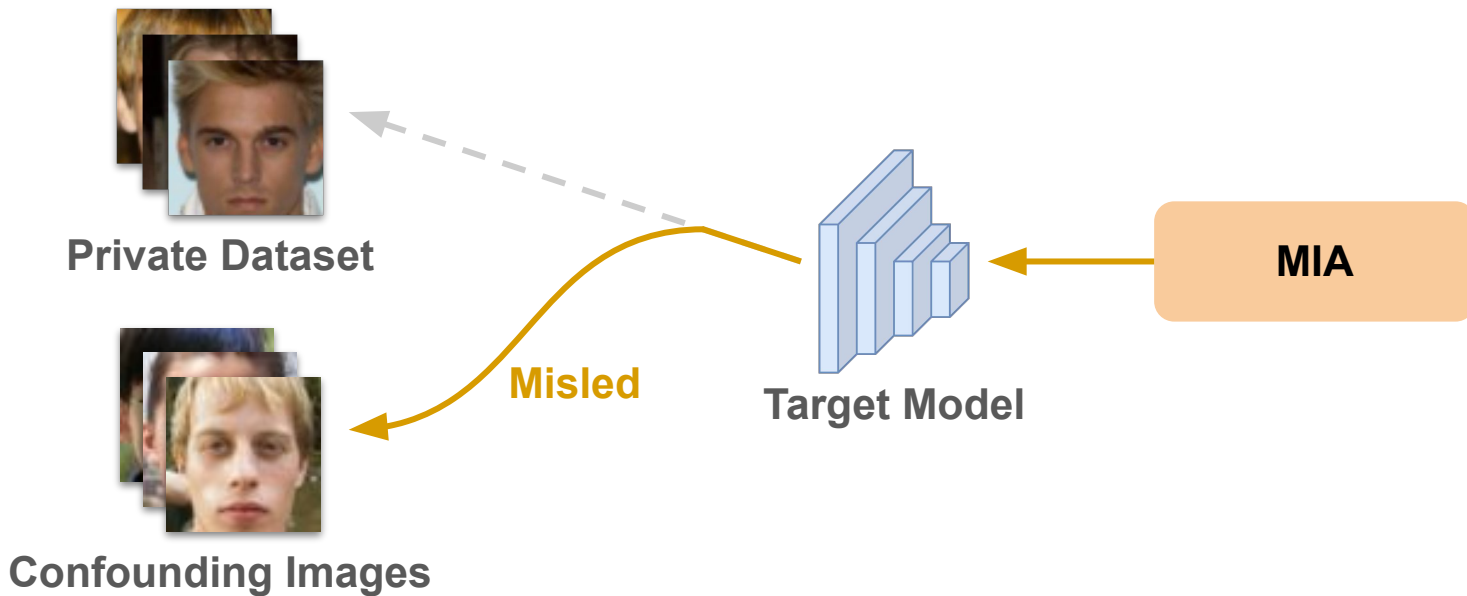
# Previous defenses are still vulnerable to recent MI attacks



> **89%**  
Attack Accuracy

# Existing trapping-based methods cannot protect all identities

NetGuard [Gong et al., 2023], DCD [Chen et al., 2023]

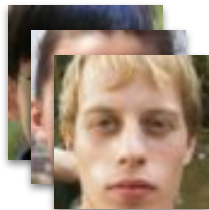


# Existing trapping-based methods cannot protect all identities

## Limitations of Existing Trapping-based Defenses

### Data Overhead

- Additional public dataset



**Confounding Images**

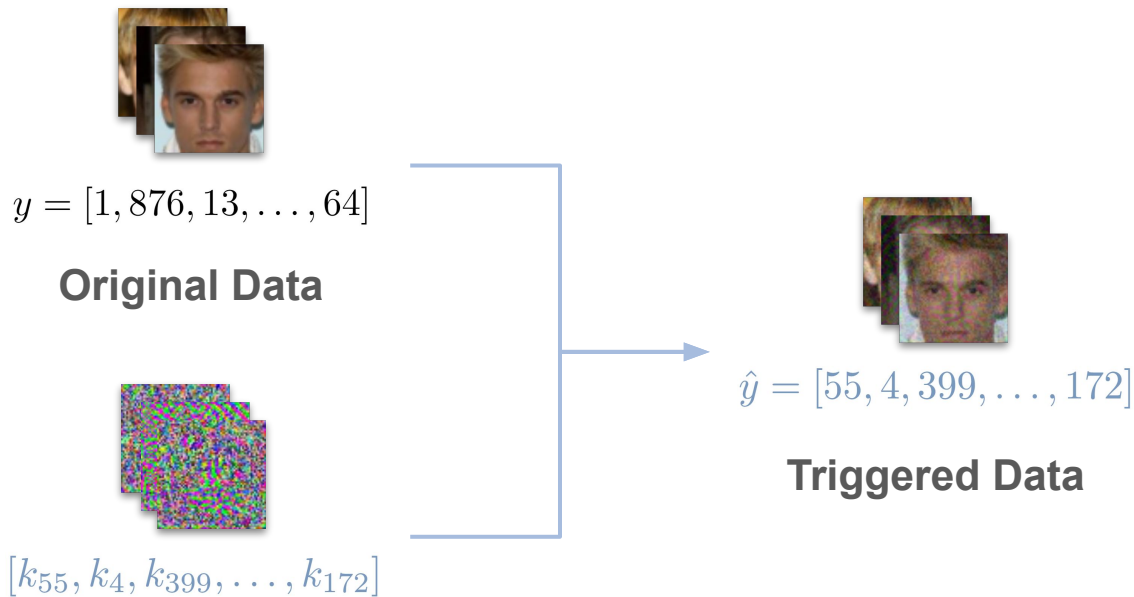
### Computational Overhead

- Training extra classifier
- Conducting shadow attacks

**Extra Classifier**

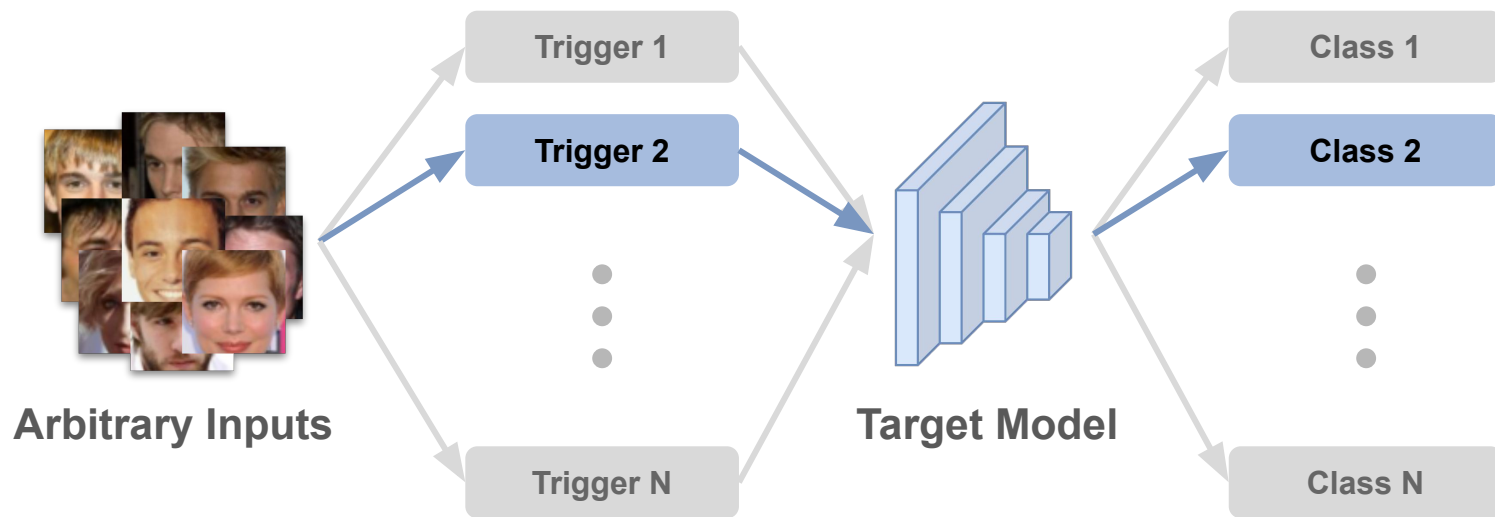
**Shadow MIA**

# Trap-MID: Mislead MI attacks by embedding trapdoors into the model



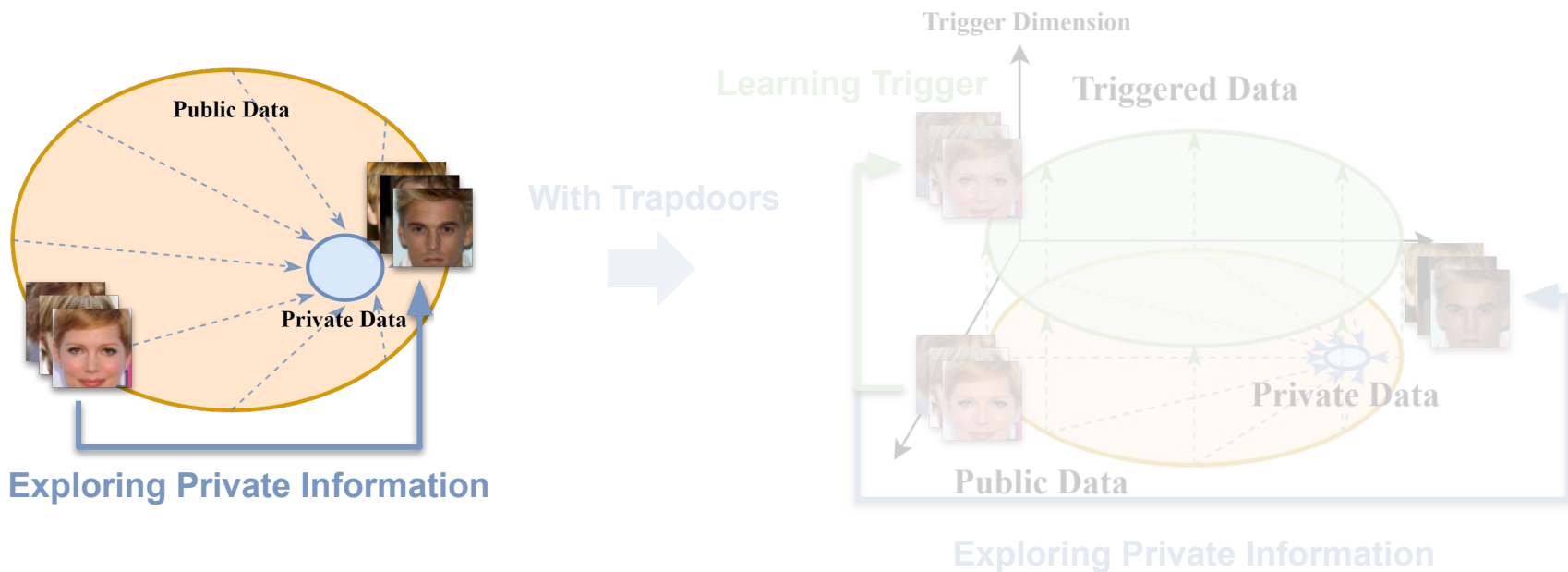
## Class-wise Trapdoor Triggers

# Trap-MID: Mislead MI attacks by embedding trapdoors into the model



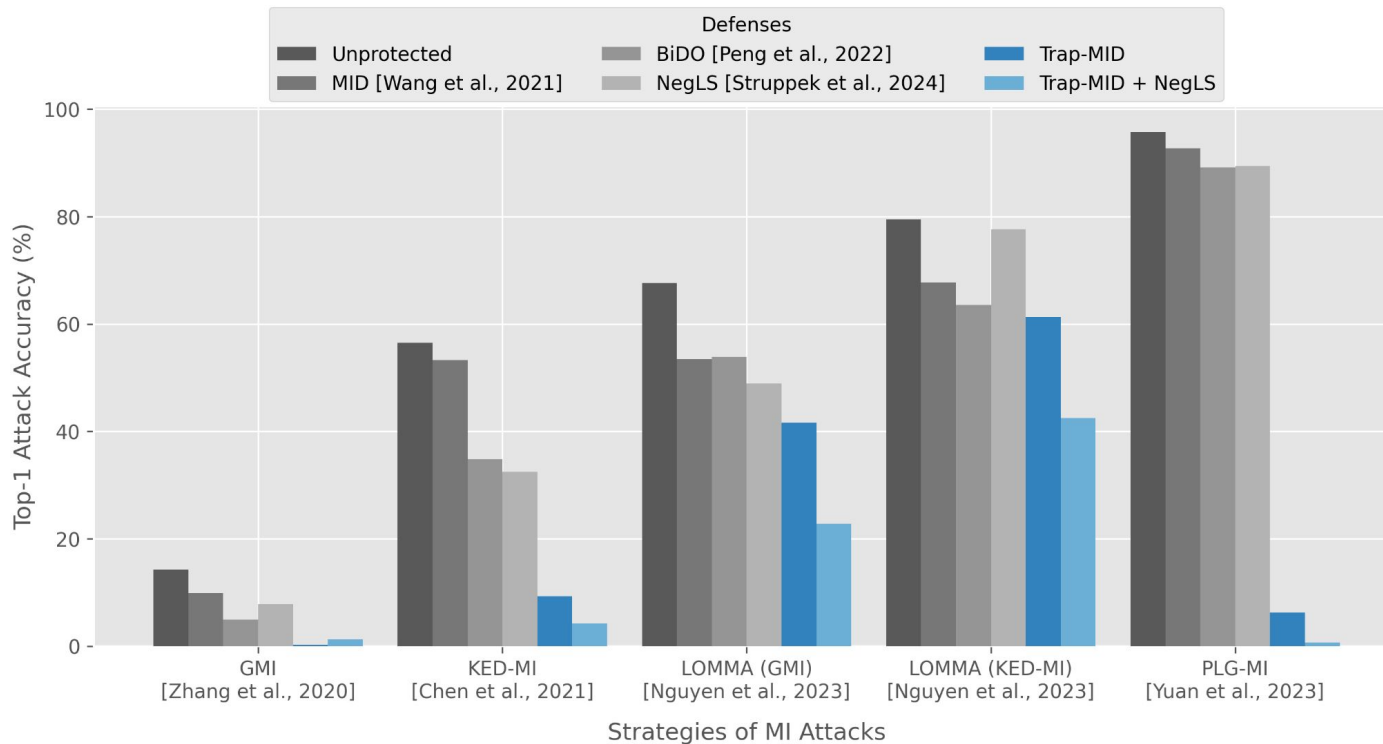


# Intuition: Trapdoor introduces shortcut for MI attacks



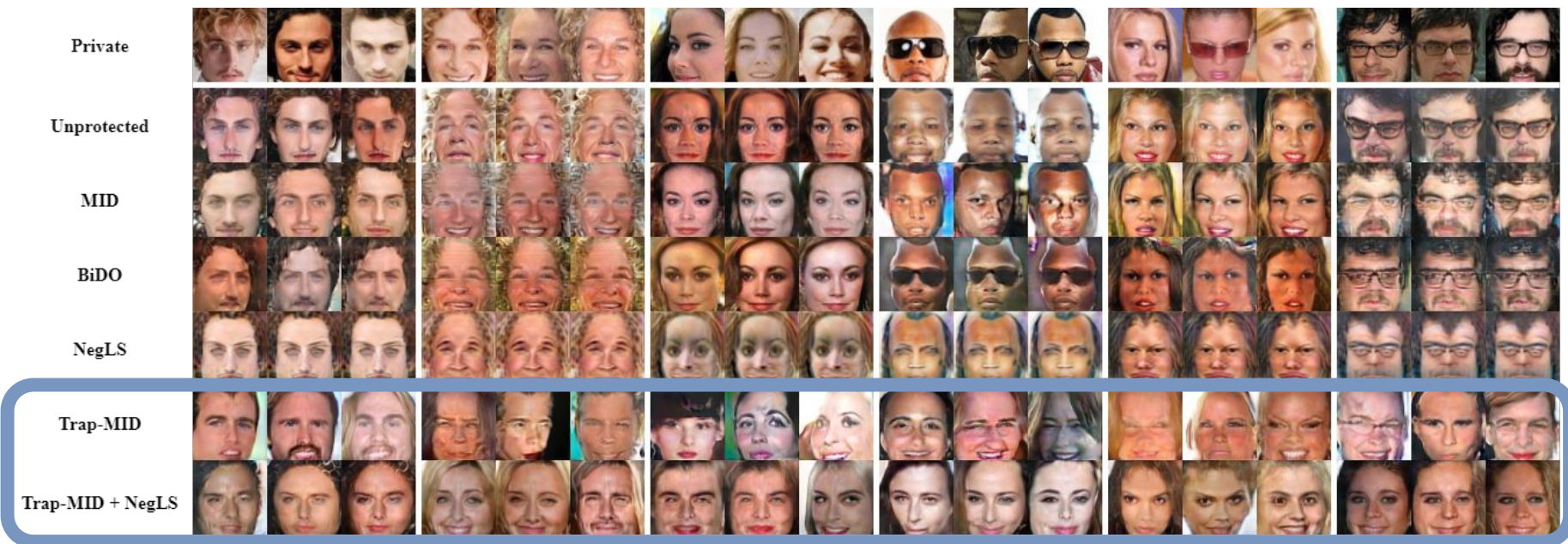


# Trap-MID provides SOTA defense against various MI attacks



## Experimental Results

# Trap-MID provides SOTA defense against various MI attacks



## Conclusion

- Trap-MID outperforms existing defense methods against MI attacks
- To the best of our knowledge, we are the first to introduce trapdoor injection technique to defend MI attacks
- Compared with existing trapping-based defense, Trap-MID preserves privacy in a more computational and data-efficient way