# A Human-in-the-Loop Fairness-Aware Model Selection Framework for Complex Fairness Objective Landscapes

**Jake Robertson**[1,2], Thorsten Schmidt[1], Frank Hutter[1,3], Noor Awad[1]

[1] University of Freiburg, [2] Zuse School ELIZA, [3] ELLIS Institute Tübingen

## In a Nutshell

Despite the recognized trade-offs among various algorithmic fairness concepts, existing fairness-aware ML methods typically focus on optimizing a single, user-specified fairness measure. This approach is problematic because:

1. Real-world FairML scenarios often involve intricate and varied stakeholder concerns, encompassing multiple fairness criteria
2. Concentrating on one fairness notion may not only compromise other pertinent fairness metrics, but also potentially result in adverse downstream effects

ManyFairHPO is a human-centered, optimization-driven framework that allows fairness practitioners to specify, rank, and optimize for multiple fairness metrics. ManyFairHPO facilitates fairness modeling decisions that effectively balance fairness objectives and reduce conflict-associated risks

## Background

| Fairness Metrics | Social Objectives |
|---|---|
| Statistical Parity (DSP) | Equality |
| Equal Opportunity (EOP) | Equity |
| Equalized Odds (EOD) | Equity |
| Inverse Distance (IND) | Individual Justice |

| Satisfied Metric \ Violated Metric | DSP | EOP | EOD | IND |
|---|---|---|---|---|
| DSP | | SFP | SFP | SFP |
| EOP | ? | | ? | ? |
| EOD | ? | ? | | ? |
| IND | ? | ? | ? | |

- Multi-objective Hyperparameter Optimization (MOHPO) involves adjusting typical ML design parameters (e.g. neural network structure) to approximate the Pareto Front of conflicting ML goals (e.g. accuracy and energy consumption)
- In fairness applications, MOHPO has been used to balance accuracy with a single, user-specified fairness criterion
- However, the established Impossibility Theorem shows that optimizing one fairness notion can unintentionally violate other relevant concepts
- This results in 1) a compromise between related social objectives and potentially 2) undesirable downstream effects (e.g. Self-Fulfilling Prophecy)

## Many-Objective Fairness-Aware Hyperparameter Optimization (ManyFairHPO)

1) **Many-Objective Optimization for a set of fairness metrics**
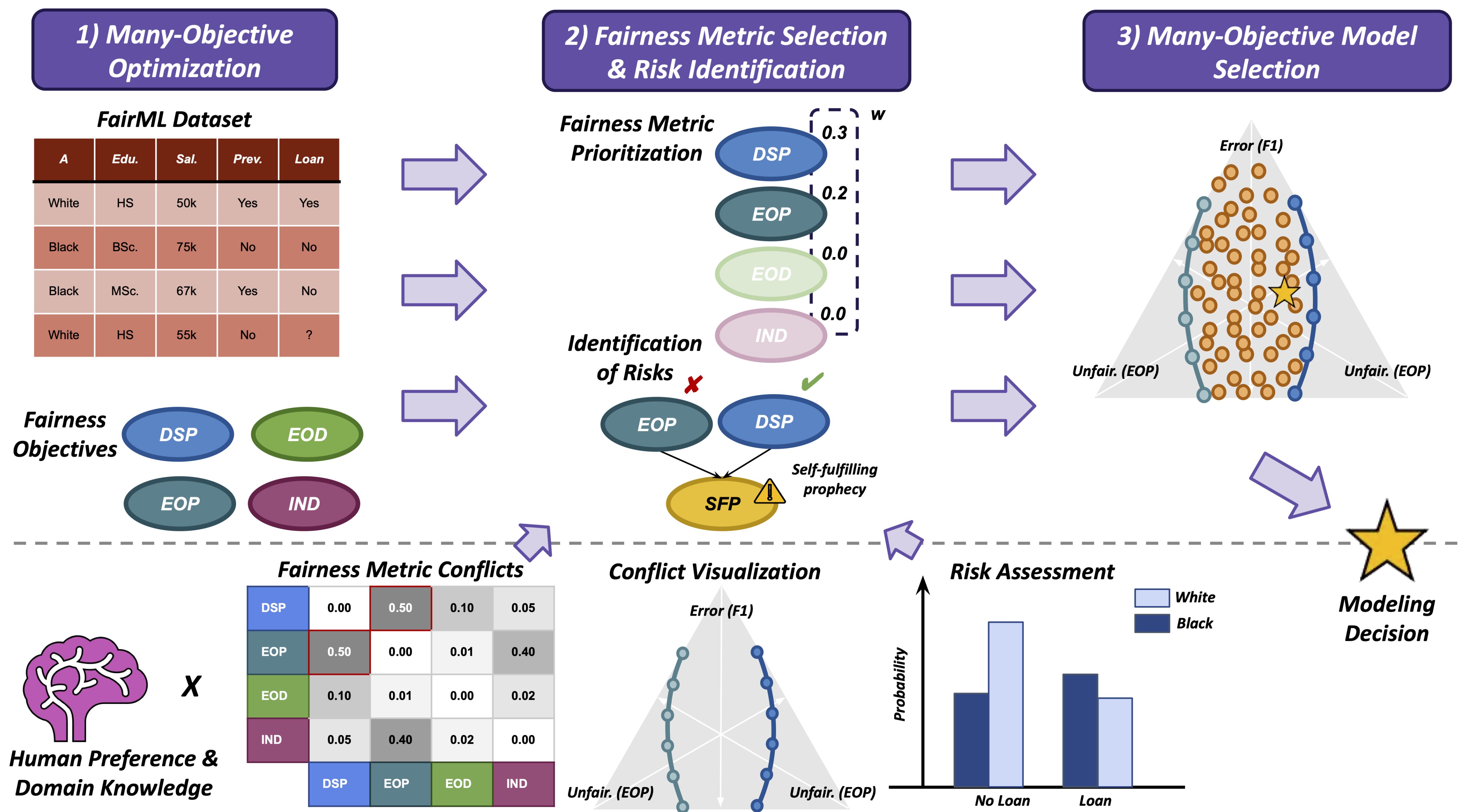   - ML dataset
   - fairness metrics

2) **Fairness Metric Selection and Risk Identification**
   - Human preferences and domain knowledge: i.e. *which set of fairness metrics are important for my task?*
   - Identified fairness metric conflicts and their associated risks: *which fairness metrics are in conflict and what might be the implications?*

3) **Many-Objective Model Selection**
   - Fairness metric weights are translated by single-objective scalarization into a single model selection decision



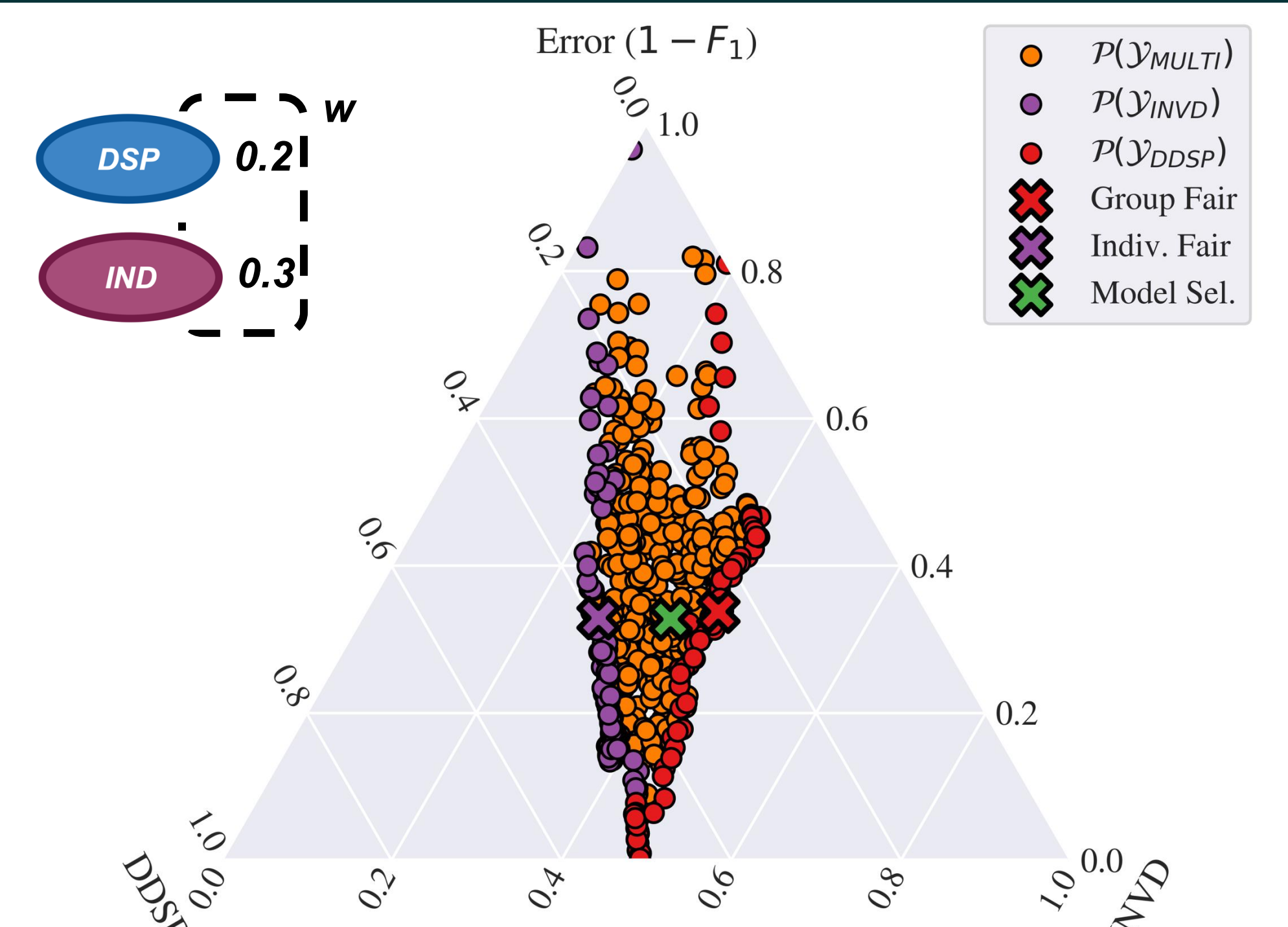## Fairness Metric Conflicts (Fairness Datasets)



**Figure 1.** Fairness metric conflicts discovered by ManyFairHPO on common fairness datasets. Problem specific conflicts can guide practitioners in selecting and prioritizing fairness metrics and identifying and assessing fairness metric conflict related risks

## Stakeholder Compromise



**Figure 2.** Given a set of fairness metric weights, $\langle 0.2, 0.3 \rangle$ for DDSP and INVD on a , ManyFairHPO selects a model (green) that balances this conflict