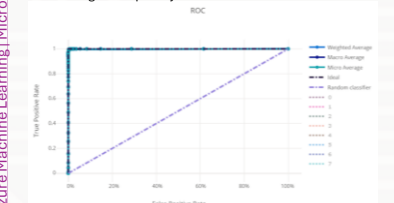
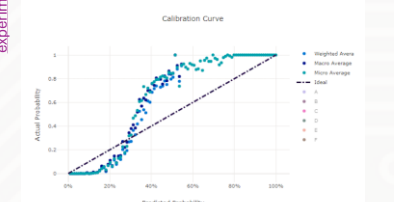


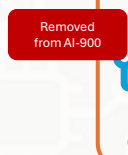
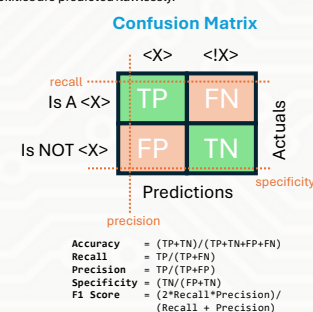
Confusion matrices visually demonstrate the systematic errors of a classification model, with "confusion" referring to the mislabeling of samples by the model.



The **receiver operating characteristic (ROC)** curve illustrates the connection between the true positive rate (TPR) and the false positive rate (FPR) as the decision threshold is varied. The area under the curve (AUC) represents the fraction of correctly classified instances.



A **calibration curve** shows a model's confidence in predictions versus the proportion of correct positive samples at each confidence level. An ideal model accurately classifies 100% of predictions with 100% confidence, 50% with 50% confidence, and 20% with 20% confidence. A perfectly calibrated model's curve matches the $y = x$ line, where probabilities are predicted flawlessly.



- Deep fakes
- Impersonation
- Harmful
- Discriminatory

1. Identify harms - what is worst case scenario?
2. Measure harms - try to generate harmful content
3. Mitigate harms - fine-tune model and filters, meta prompts, grounding data
4. Operate responsibly - monitor, feedback, track telemetry