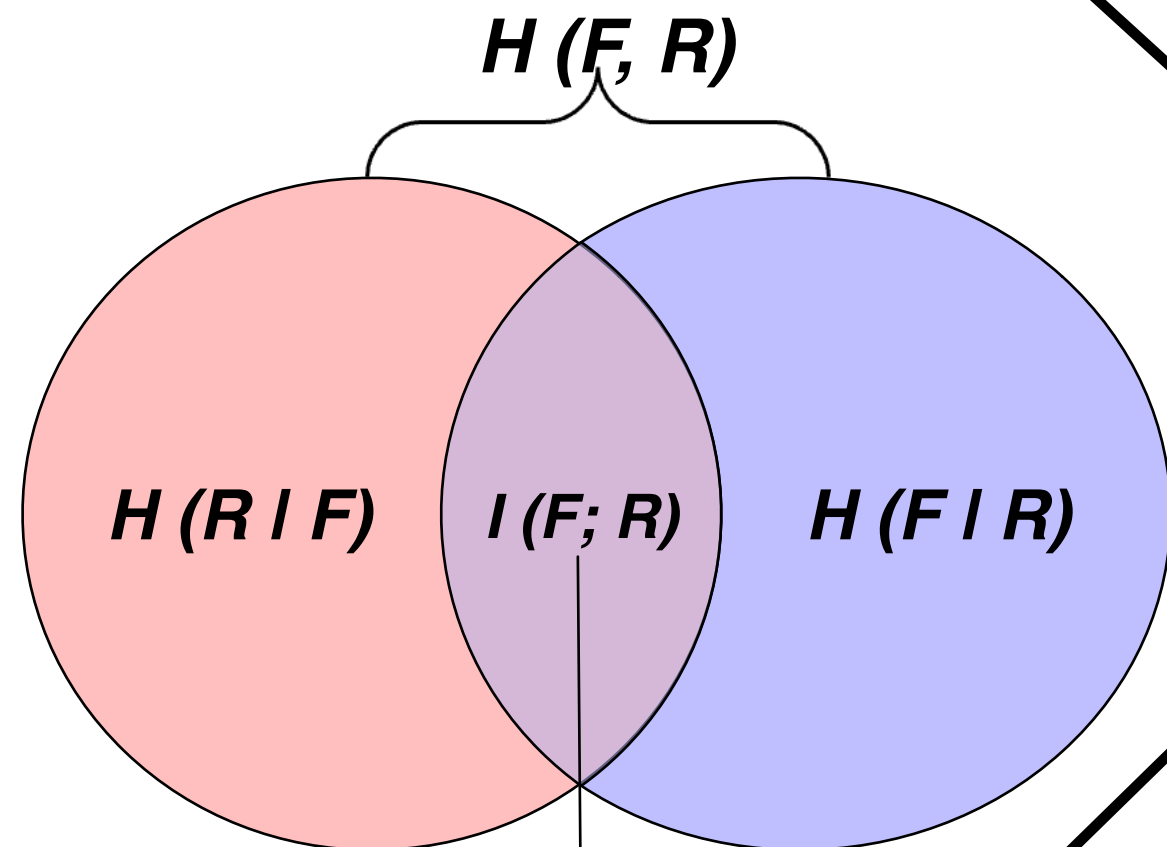


Retain
Dataset

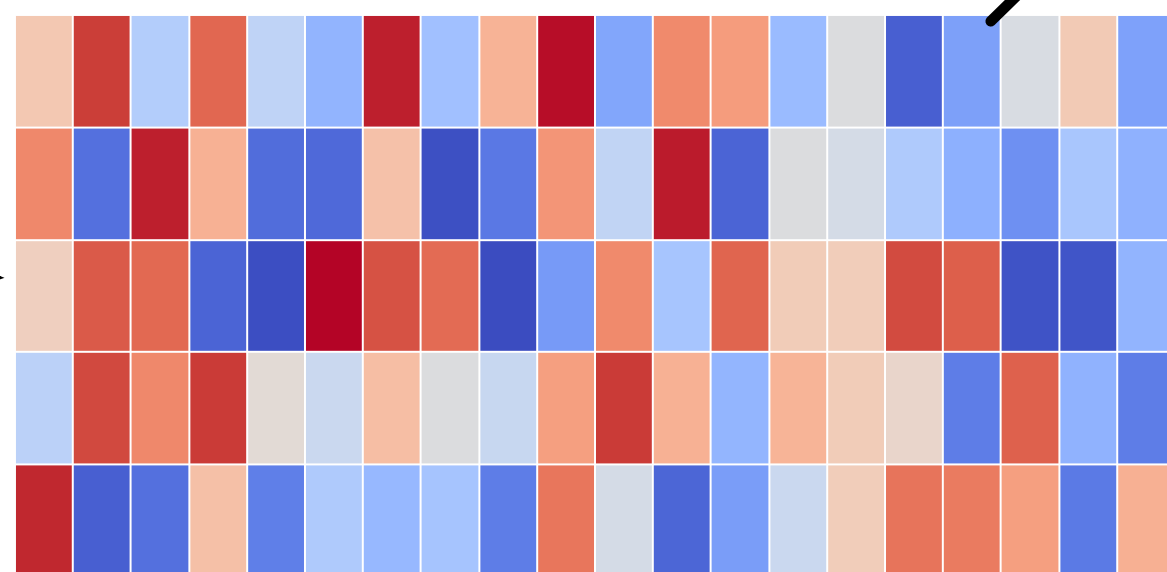


Retain Activation (R)



Mutual Information

Forget
Dataset

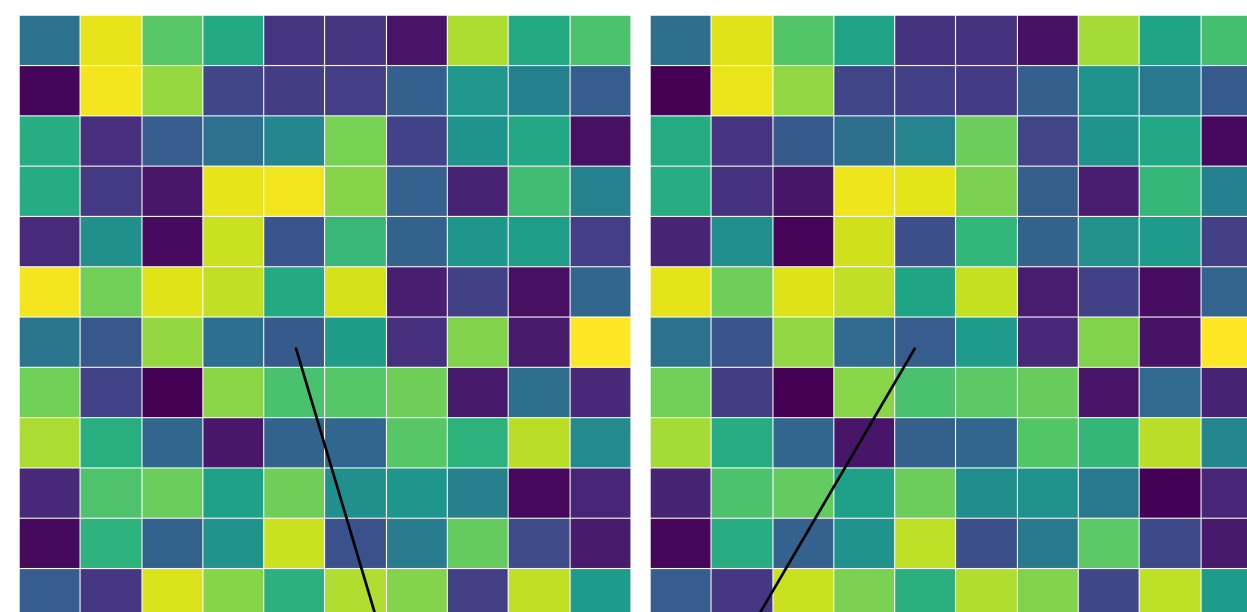


Forget Activation (F)

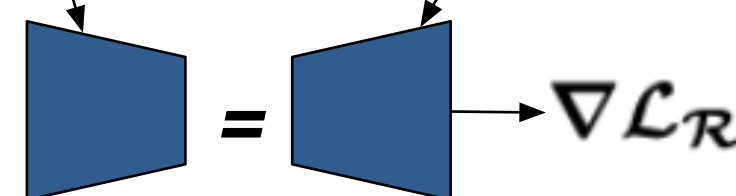
Step 1: Identify Low-Entanglement Parameters via Information-Theoretic Guidance

Updated R

Frozen R

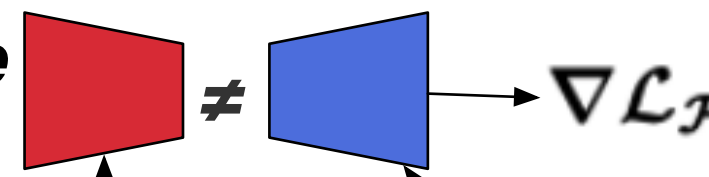


Retain
LOSS



Step 2.1: Contrastive Representation Unlearning

Forgettable
LOSS

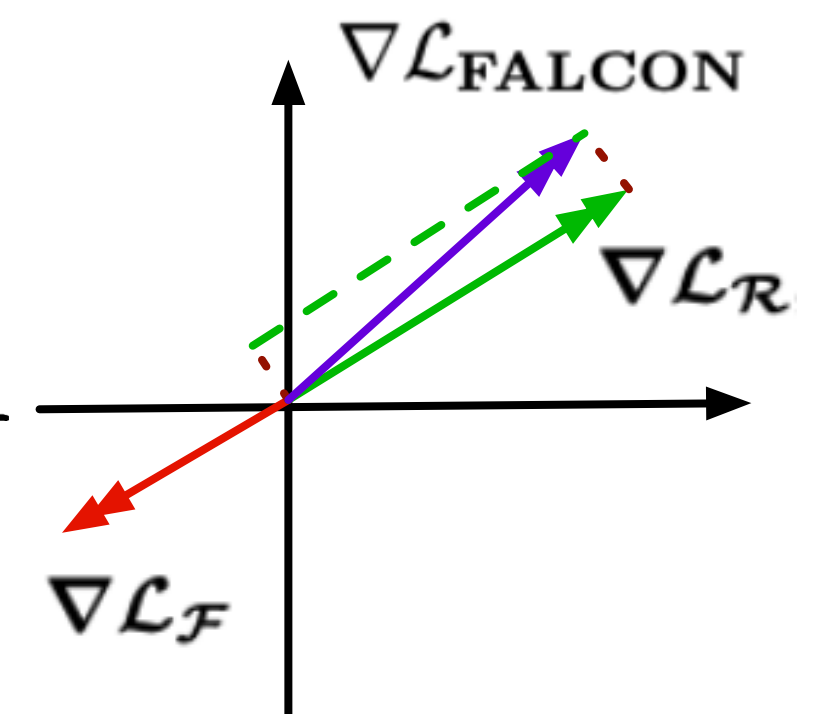
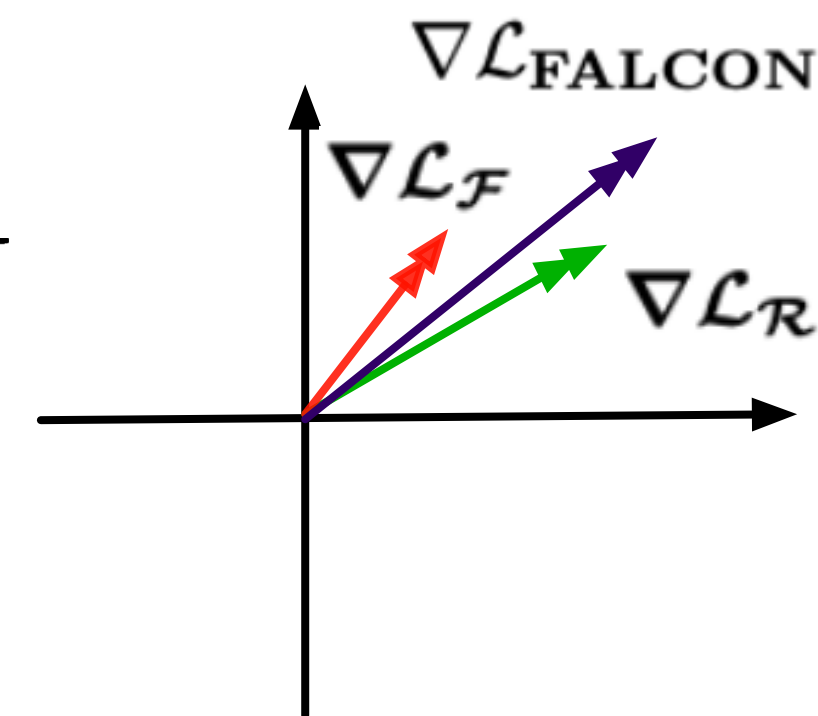


$\text{COS}(\cdot) > 0$

Yes

No

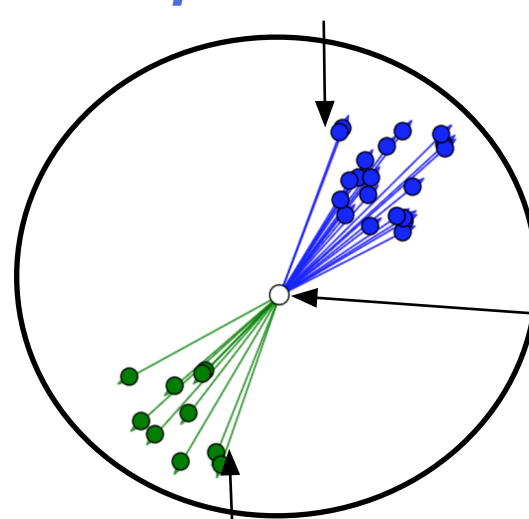
Step 2.2: Orthogonalizing Gradient for Conflict Resolution



Step 3: Model Unlearning



Unwanted
Representation



Steer Unwanted
Representation Away via
Principal Offset Vector

Updated F

Frozen F

