

# Human-Driven FOL Explanations of Deep Learning

**Gabriele Ciravegna,**

SAILAB, University of Siena, Italy



**SAILab**  
Siena Artificial Intelligence Lab



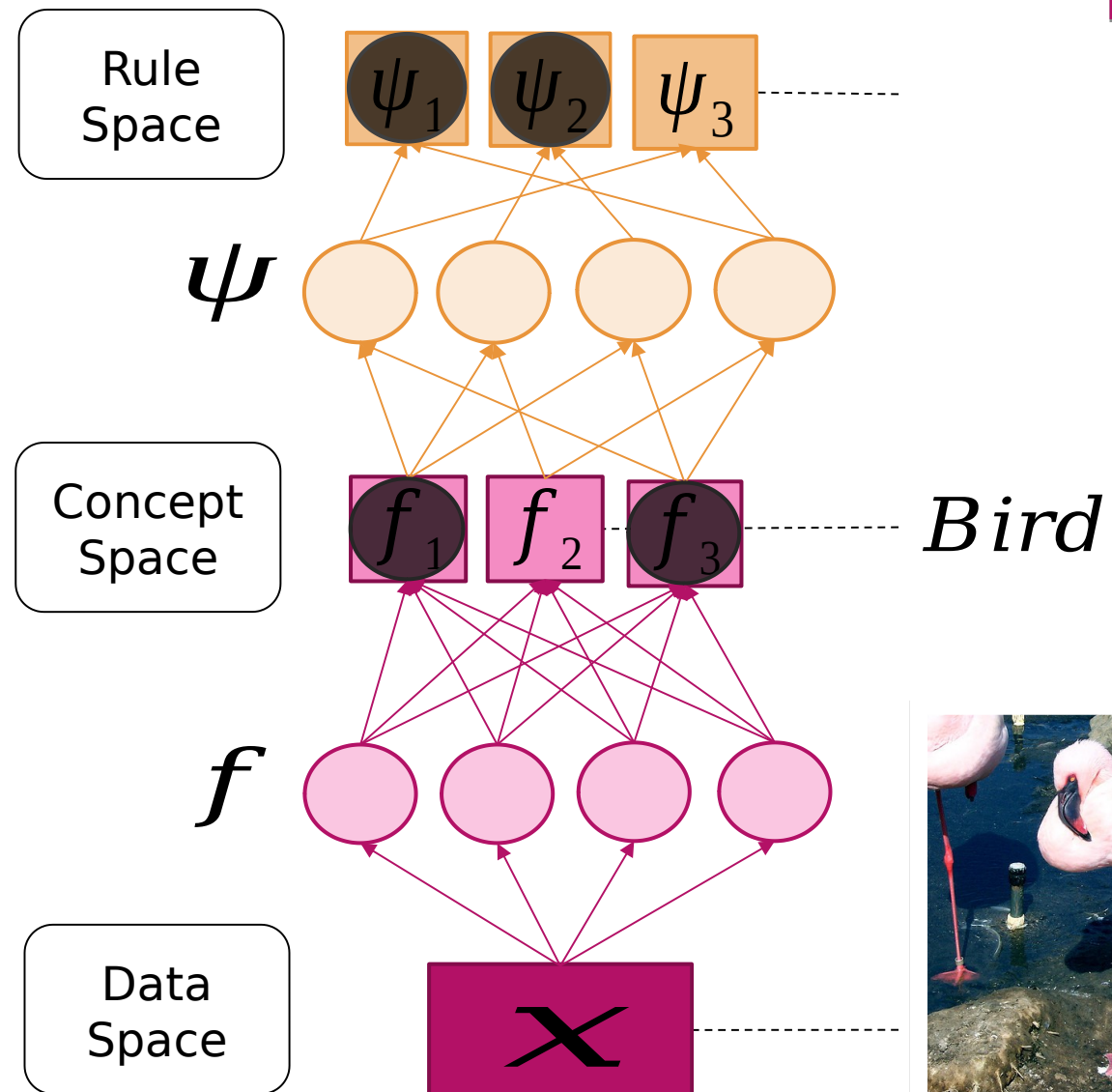
**UNIVERSITÀ  
DI SIENA 1240**



NEURAL NETWORKS HAVE BECOME INCREDIBLY  
POWERFUL

# Network explaining network

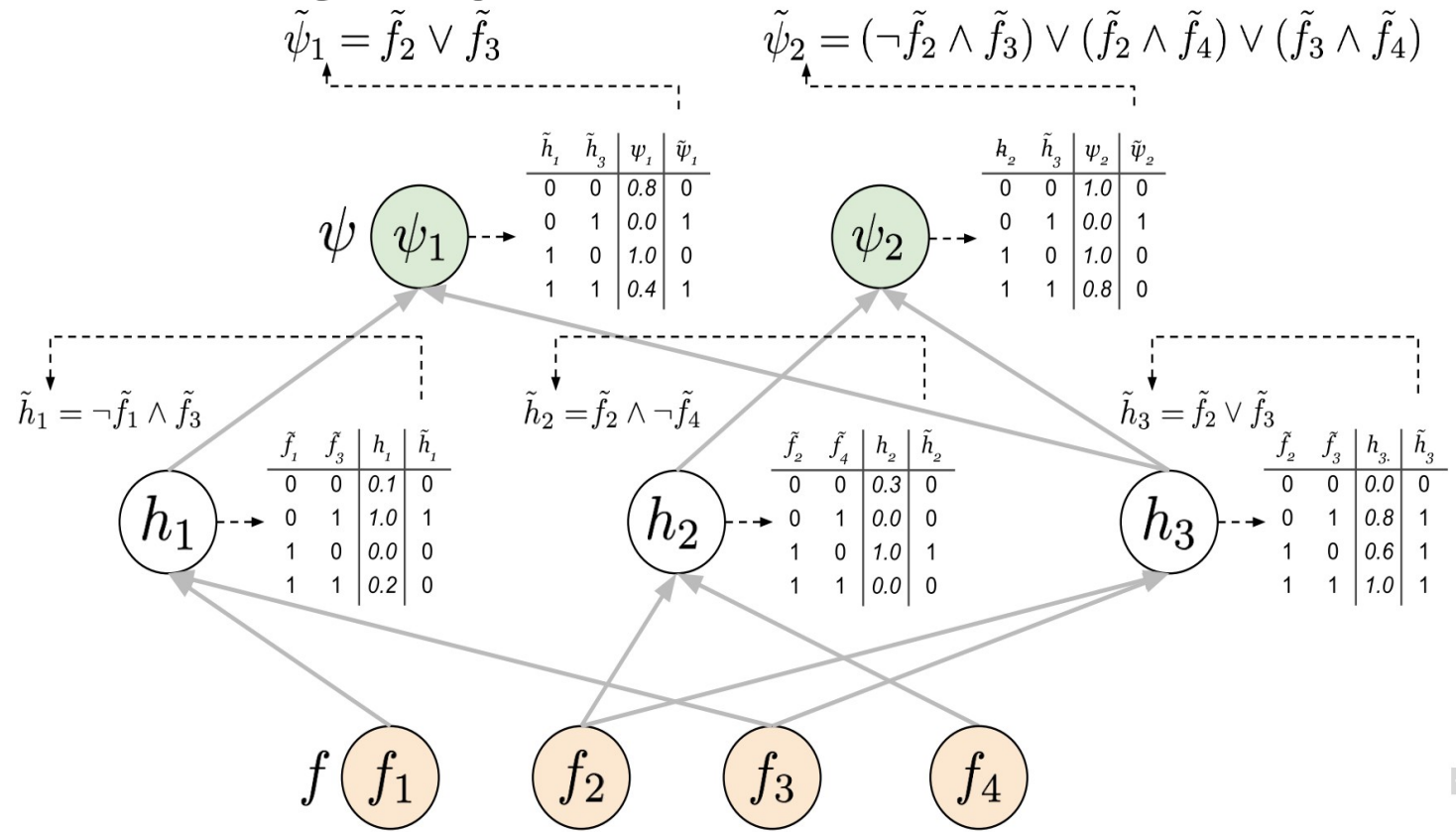
- Provides example-based explanation
- Extract new knowledge



# Network Pruning

- ▶ Strong L1 regularization
- ▶ Progressive pruning of the lowest weights
- ▶ **Low Fan-In Neurons**

## Explainable Network



# How does it work?

$$f^*, \psi^* = \arg \min_{f, \psi} \{ \mathcal{U}(f) + D(\psi, f) \}$$

$$\mathcal{U}(f) = \left\{ \sum_{j=1}^c \sum_{\mathbf{x}_k \in X_{\phi_j}} \hat{\phi}_j(f(\mathbf{x}_k)) + \gamma_f \|f\| \right\}$$

$$\hat{\phi}_j(f(\mathbf{x})) = \|f(\mathbf{x}) - \mathbf{y}(\mathbf{x})\|^2$$

$$D(\psi, f) = \{ -\hat{I}_{Y, \psi}(\psi, f, X) + \gamma_{\psi} \|\psi\| \}$$

$$I_{Y, \psi}(\psi, f) = H_{\psi}(\psi, f) - H_{\psi|Y}(\psi, f)$$



# Regularization Effects

**Baseline**

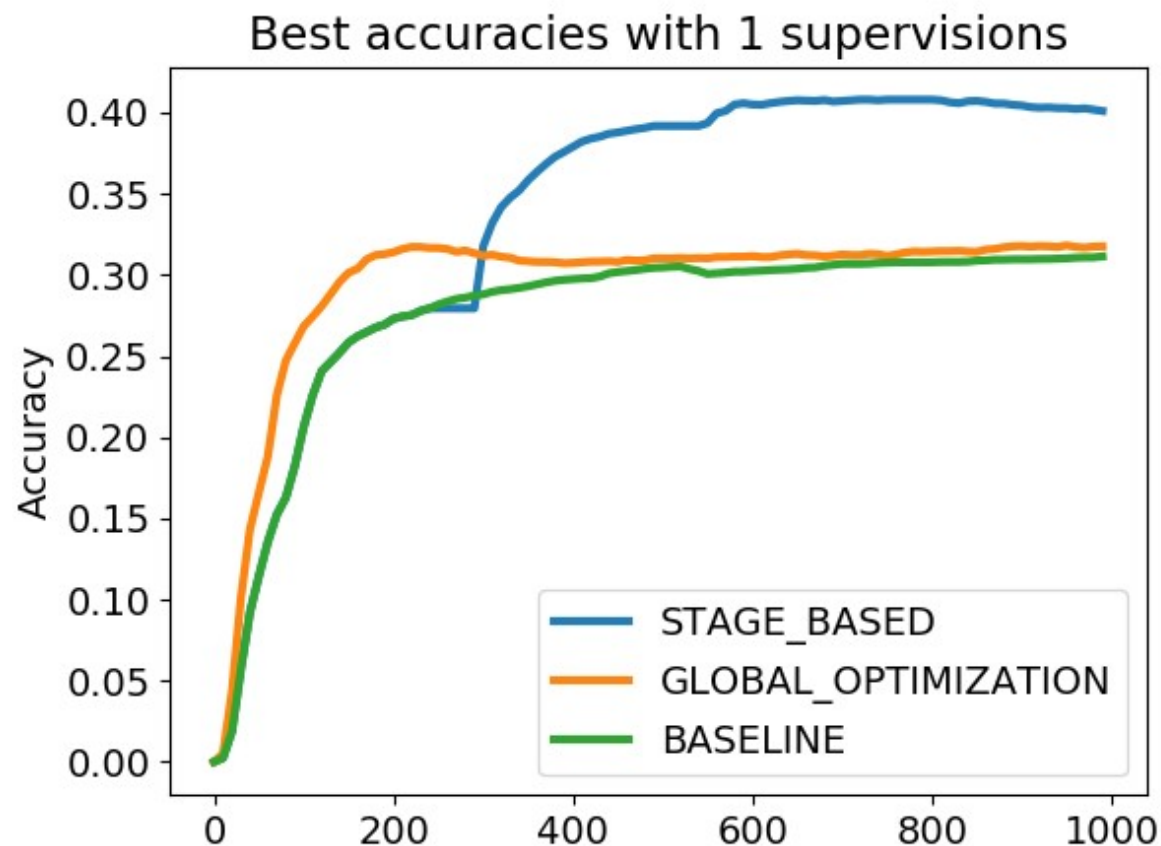
$$U(F)$$

**Global  
Optimization**

$$U(F) + D(f, \psi)$$

**Stage-based  
Optimization**

# Performance Improvements - MNIST



# Overall Results - MNIST

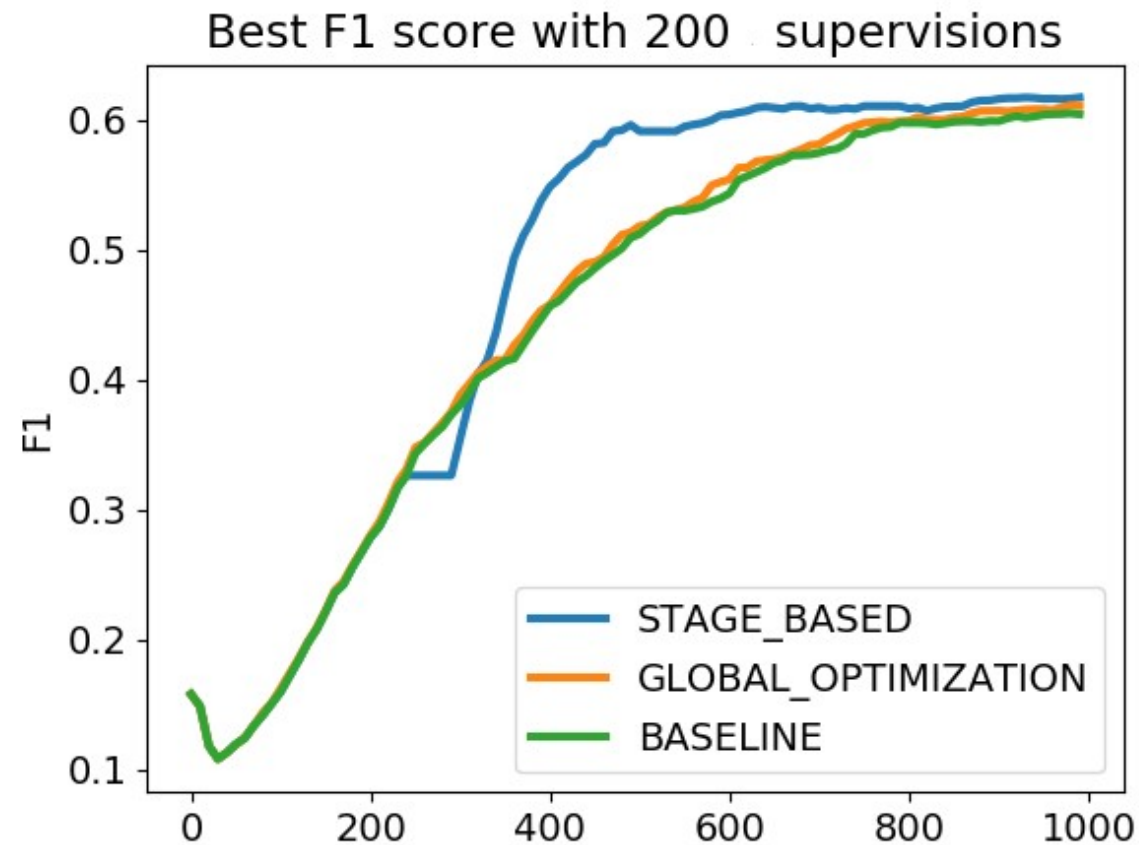
# Labeled Examples	Baseline	Stage-based Optimization	Global Optimization
1	31.7 $\pm$ 2.1 %	<b>46.5 <math>\pm</math> 1.4 %</b>	38.7 $\pm$ 3.7%
10	68.6 $\pm$ 0.4 %	71.5 $\pm$ 0.9%	<b>71.6 <math>\pm</math> 1.7%</b>
100	84.5 $\pm$ 0.4 %	<b>86.7 <math>\pm</math> 0.3 %</b>	86.0 $\pm$ 0.1%

## Learned Rules

(a)	Odd $\dot{\vee}$ Even Nine $\wedge \neg$ Even	Eight $\dot{\vee}$ Zero Six $\wedge \neg$ Odd
(b)	Odd $\wedge \neg$ Even $\wedge$ [(One $\dot{\vee}$ Five $\dot{\vee}$ Seven $\dot{\vee}$ Nine) $\vee$ ( $\neg$ One $\wedge \neg$ Five $\wedge \neg$ Seven $\wedge \neg$ Nine)] Even $\wedge \neg$ Odd $\wedge \neg$ One $\wedge \neg$ Five $\wedge \neg$ Nine $\wedge$ (Eight $\vee \neg$ Eight)	



# Performance Improvements - PascalPart



# Overall Results - PascalPart

# Labeled Examples	Baseline	Stage-based Optimization	Global Optimization
10	54.0 $\pm$ 0.3%	<b>56.8 <math>\pm</math> 0.2%</b>	55.8 $\pm$ 0.7%
50	60.5 $\pm$ 0.2%	<b>62.0 <math>\pm</math> 0.4%</b>	61.3 $\pm$ 0.5%
Whole dataset	62.6 $\pm$ 0.3%	<b>63.8 <math>\pm</math> 0.2%</b>	63.7 $\pm$ 0.4%

## Learned Rules

(a) Hand $\wedge$ Foot $\wedge$ $\neg$ Beak, Handlebar $\wedge$ $\neg$ Boat,	$\neg$ Foot $\wedge$ Dog, $\neg$ AirplaneBody $\wedge$ Cow
(b) Beak $\wedge$ $\neg$ Ear $\wedge$ $\neg$ Nose, TvMonitor $\wedge$ (Screen $\vee$ Table)	(Cat $\vee$ Dog) $\wedge$ ( $\neg$ Foot $\vee$ Paw), Aeroplane $\wedge$ (Engine $\vee$ Stern)

# Following steps

## Improve rules quality

- From Local to Global rules:
- Predicate explaining rules:

## Framework Applications

- Adversarial Defense:  
Learned constraints as Attack  
detector

# From Local to Global Explanations

$$\forall x \in X_j, \hat{\psi}_j(f(x)) \quad X = \bigcup_{j=1}^m X_j$$



**CNF  
Conversion**

$$\hat{\Psi} = \bigvee_{j=1}^m \hat{\psi}_j(f(x)) \equiv \bigwedge_{k=1}^K \hat{\psi}'_k(f(x))$$

$$\forall x \in X, \hat{\psi}'_k(f(x))$$

# From Local to Class-Driven Explanations

**Local Rule  
rewrite**  $\hat{\psi}_i(f(x)) = 1_{\psi_i}(f(x)) \leftrightarrow \psi_i(f(x)), \quad \forall x \in X,$

**Enforce Support**  $1_{\hat{\psi}_i(f(x))} = 1_{\hat{f}_i(x)}$



**Class-Driven  
Rule**  $\hat{\psi}_i(f(x)) = \hat{f}_i(x) \leftrightarrow \psi_i(f_k(x)), \quad \forall x \in X, \quad i, k \in (1, \dots, n), \quad k \neq i$

$\hat{\psi}_i(f(x)) = Man(x) \leftrightarrow Head(x) \vee Hand(x) \vee Foot(x), \quad \forall x \in X$

# IFF $\leftrightarrow$ or IF $\rightarrow$ ?

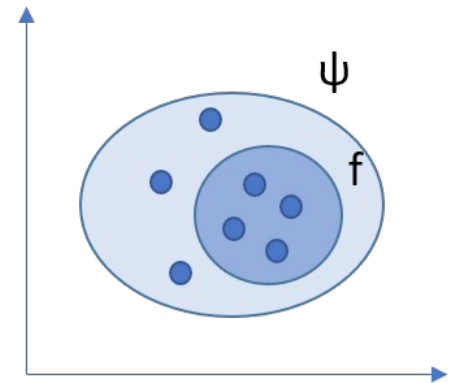
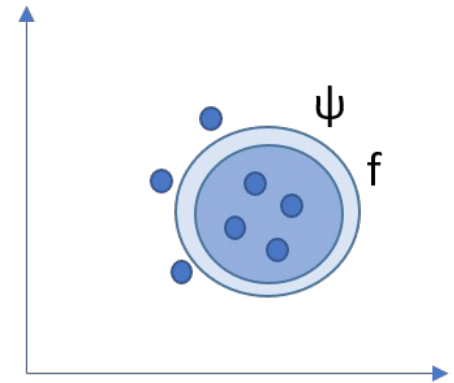
- **IFF  $\leftrightarrow$  rules:**

$$1_{\hat{\psi}_i(f(x))} = 1_{\hat{f}_i(x)}$$

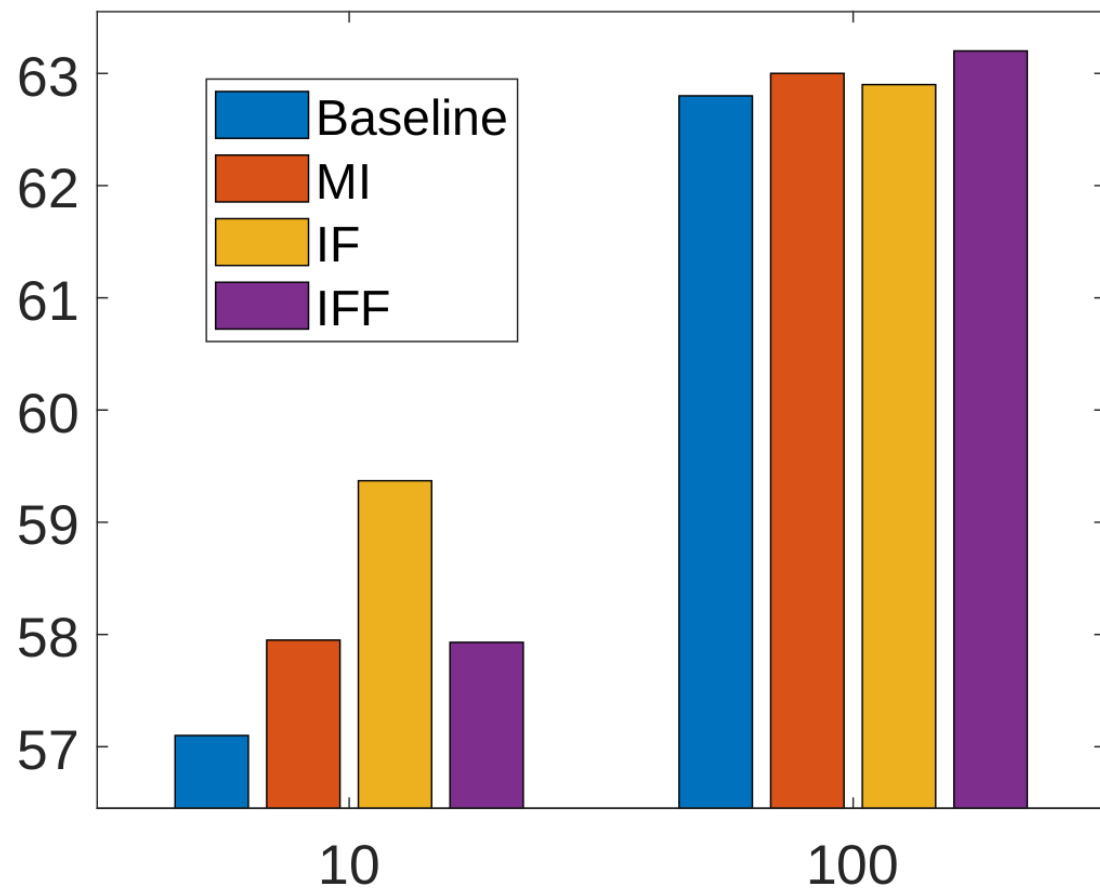
- **IF  $\rightarrow$  rule:** 
$$L_{\leftrightarrow}(\psi, f, X) = \sum_{i \in S, x \in X} |f_i(x) - \psi_{h(i)}(f(x))|$$

$$1_{\hat{\psi}_i(f(x))} \supseteq 1_{\hat{f}_i(x)}$$

$$L_{\rightarrow}(\psi, f, X) = \sum_{i \in P, x \in X} \max\{0, f_i(x) - \psi_{h(i)}(f(x))\}$$





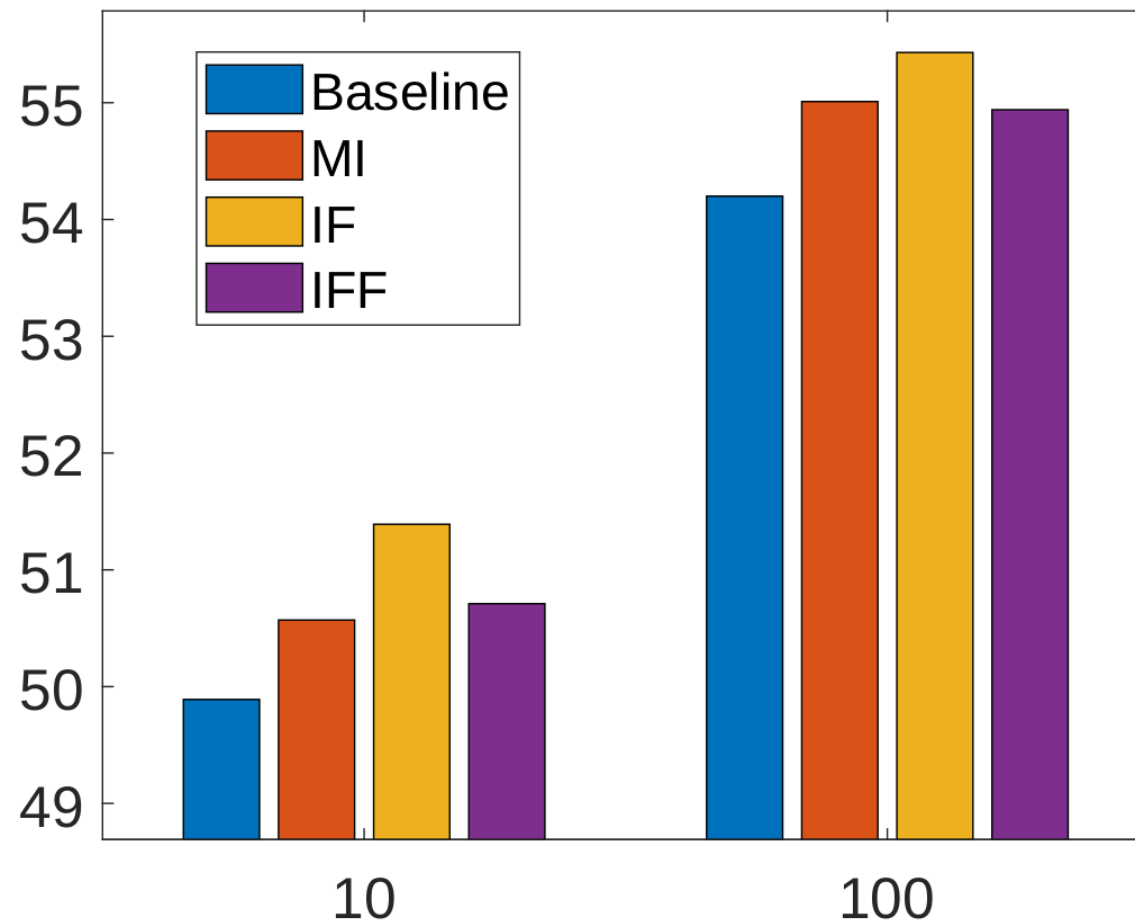


PascalPart  
Performances

LOCAL	$\forall x \in X_{\psi_i}, \text{Beak} \vee \text{Bird}$ $\forall x \in X_{\psi_j}, \text{Headlight} \vee \text{Plate}$ $\forall x \in X_{\psi_k}, \text{Cat} \vee \text{Horse}$
GLOBAL	$\forall x, \text{AeroplaneBody} \vee \text{Beak} \vee \text{Bird} \vee \text{Table} \vee \text{Plant}$ $\vee \text{Car} \vee \text{Headlight} \vee \text{Motorbike} \vee \text{Muzzle} \vee \text{Train}$ $\vee \text{Chainwheel} \vee \neg \text{Aeroplane}$ $\forall x, \neg \text{Horse} \vee \text{AeroplaneBody} \vee \text{Beak} \vee \text{Bird} \vee \text{Train}$ $\vee \text{Car} \vee \text{Chainwheel} \vee \text{Headlight} \vee \text{Muzzle} \vee \text{Table}$ $\vee \text{Motorbike} \vee \text{Plant}$
CLASS- DRIVEN $IF \rightarrow$	$\forall x, \text{Car} \rightarrow \text{Backside} \vee \text{Mirror} \vee (\text{Window} \wedge \neg \text{Coach})$ $\forall x, \text{Bicycle} \rightarrow \text{Saddle} \vee \text{Handlebar}$ $\forall x, \text{Train} \rightarrow \text{Coach} \vee \text{TrainHead}$
CLASS- DRIVEN $IFF \leftrightarrow$	$\forall x, \text{Horse} \leftrightarrow (\text{Hoof} \wedge \text{Ear}) \vee (\text{Hoof} \wedge \text{Neck})$ $\forall x, \text{Bird} \leftrightarrow \text{Beak} \wedge \neg \text{Horn}$ $\forall x, \text{Bicycle} \leftrightarrow (\text{Chainwheel} \wedge \neg \text{Cow} \wedge \text{Handlebar})$ $\vee (\text{Chainwheel} \wedge \neg \text{Cow} \wedge \text{Saddle})$

# Rules PascalPart

# CelebA Performances



# Rules CelebA

LOCAL	$\forall x \in X_{\psi_i}, \text{Bangs} \wedge \neg \text{Bald}$ $\forall x \in X_{\psi_j}, \text{StraightHair} \wedge \text{BushyEyebrows}$ $\forall x \in X_{\psi_k}, \text{Female} \wedge \text{Attractive}$
GLOBAL	$\forall x, \text{Bangs} \vee \text{BlondHair} \vee \text{Blurry} \vee \text{Goatee}$ $\vee \text{StraightHair} \vee \text{WearHat}$ $\vee \neg \text{Attractive} \vee \neg \text{Female} \vee \neg \text{Male}$ $\forall x, \text{Blurry} \vee \text{Goatee} \vee \text{WearHat}$ $\vee \neg \text{Attractive} \vee \neg \text{BlackHair} \vee \neg \text{BlondHair}$ $\vee \neg \text{Female} \vee \neg \text{StraightHair}$
CLASS- DRIVEN $IF \rightarrow$	$\forall x, \text{Attractive} \rightarrow \text{PaleSkin} \vee \text{RosyCheeks}$ $\vee (\neg \text{Blurry} \wedge \neg \text{Chubby})$ $\forall x, \text{Beard} \rightarrow \text{Goatee} \vee \text{Sideburns}$ $\forall x, \text{Old} \rightarrow \text{GrayHair} \vee \neg \text{Attractive}$
CLASS- DRIVEN $IFF \leftrightarrow$	$\forall x, \text{Bald} \leftrightarrow \neg \text{BlackHair} \wedge \neg \text{BrownHair}$ $\wedge \neg \text{StraightHair} \wedge \neg \text{WavyHair}$ $\forall x, \text{NotBald} \leftrightarrow \text{Bangs} \vee \text{BrownHair} \vee \text{WavyHair}$ $\forall x, \text{Male} \leftrightarrow \neg \text{WearLipstick} \wedge \neg \text{WearNecklace}$





SAILab  
Siena Artificial Intelligence Lab

Thank you for the  
attention

---