

# Differentiable Logics for Machine Learning: What difference do they make?

Thomas Flinkow<sup>1</sup>, Barak A. Pearlmuter<sup>1,2</sup>, Rosemary Monahan<sup>1,2</sup>

<sup>1</sup>Department of Computer Science, Maynooth University

<sup>2</sup>Hamilton Institute, Maynooth University

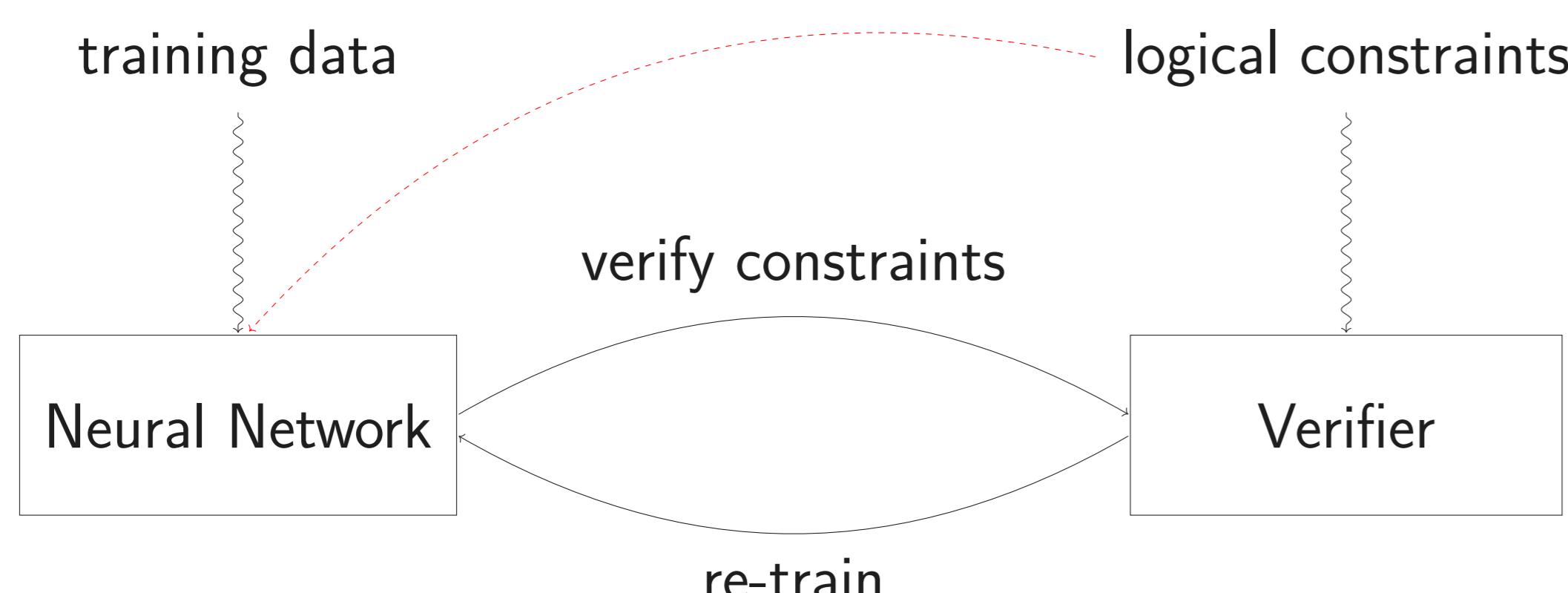
thomas.flinkow@mu.ie

## 1. Motivation: Correct-by-construction ML models

**Task:** train a neural network  $\mathcal{N}$  to satisfy constraint  $\phi$ .

**Train:** given data, labels, and loss function, iteratively update weights.

**Verify:**  $\alpha, \beta$ -CROWN, Marabou, NNV, ERAN, ...



## 2. Background: Property-driven ML

Given data  $x_0$ , label  $y$ , and constraint  $\phi$ , obtain optimal weights  $\theta^+$  by

$$\theta^+ = \arg \min_{\theta} \alpha \mathcal{L}_{CE}(x_0, y) + \beta \mathcal{L}_C(x_0, y, \phi).$$

DL2: approximate  $\forall x. x \models \phi$  by finding  $x^*$  such that  $x^* \not\models \phi$ .

- Approximate counterexample *outside* of training set using PGD:

$$x^* = \arg \max_{x \in \|x - x_0\|_\infty \leq \epsilon} \mathcal{L}_C(x_0, x, y, \phi)$$

- Use this counterexample in training:

$$\theta^+ = \arg \min_{\theta} \alpha \mathcal{L}_{CE}(x_0, y) + \beta \mathcal{L}_C(x_0, x^*, y, \phi).$$

## 3. Background: Differentiable Logics

- DL2:  $\llbracket \cdot \rrbracket_{DL2} : \Phi \rightarrow [0, \infty)$ , where  $\llbracket \top \rrbracket_{DL2} = 0$ , and

$$\llbracket x \leq y \rrbracket_{DL2} := \max\{x - y, 0\}$$

$$\llbracket \phi \wedge \psi \rrbracket_{DL2} := \llbracket \phi \rrbracket_{DL2} + \llbracket \psi \rrbracket_{DL2}$$

$$\llbracket \phi \vee \psi \rrbracket_{DL2} := \llbracket \phi \rrbracket_{DL2} \cdot \llbracket \psi \rrbracket_{DL2}.$$

- Fuzzy Logics:  $\llbracket \cdot \rrbracket_L : \Phi \rightarrow [0, 1]$ , where  $\llbracket \top \rrbracket_L = 1$  and  $\llbracket \perp \rrbracket_L = 0$

Logic	Conjunction	Disjunction	Implication
Gödel	$\min\{x, y\}$	$\max\{x, y\}$	$\begin{cases} 1, & \text{if } x < y, \\ y, & \text{else.} \end{cases}$
Kleene-Dienes			
Łukasiewicz	$\max\{0, x + y - 1\}$	$\min\{1, x + y\}$	$S(N(x), y)$
Reichenbach			
Goguen	$xy$	$x + y - xy$	$\begin{cases} 1, & \text{if } x < y, \\ y^x, & \text{else.} \end{cases}$

## 4. Methods

General-purpose framework implementation in PyTorch:

```

def train(..):
    for _, (inputs, labels) in enumerate(train_loader):
        outputs = NN(inputs)
        ce_loss = F.cross_entropy(outputs, labels)

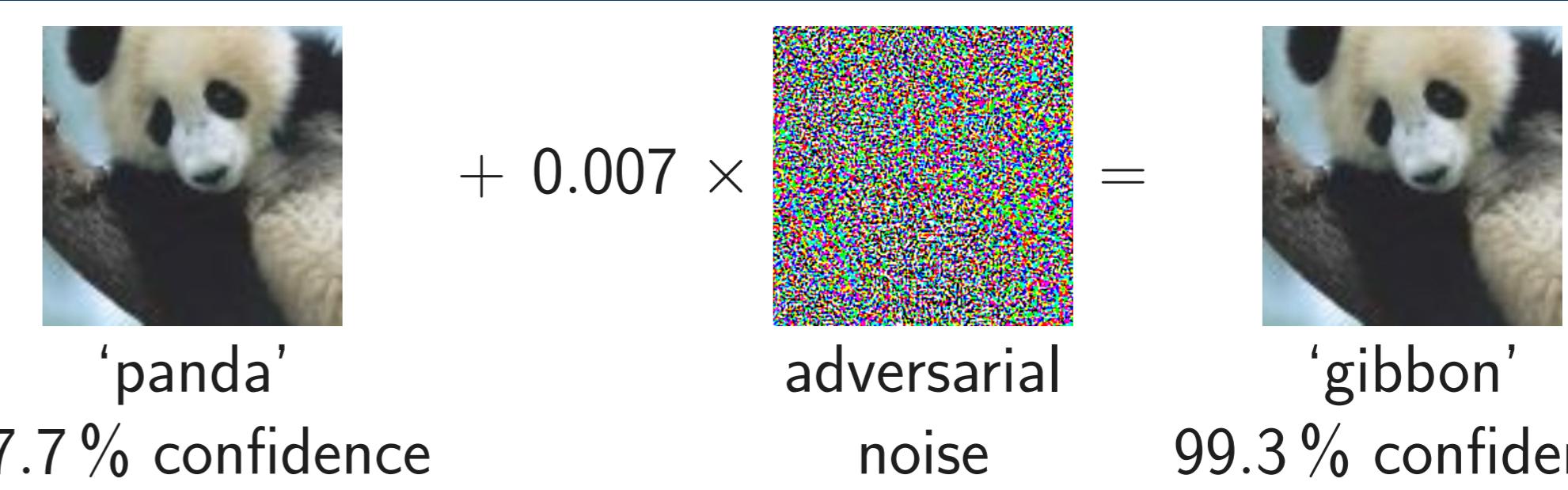
        adv = pgd.attack(NN, inputs, labels, constraint)
        dl_loss = constraint.eval(NN, inputs, adv, labels)

        loss = alpha * ce_loss + beta * dl_loss

        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
  
```

(<https://github.com/tflinkow/comparing-differentiable-logics>)

## 5. Experiment – Local Robustness Constraint



**Constraint:** A neural network is *locally robust* in input  $x_0$ , if

$$\forall x. \|x - x_0\|_\infty \leq \varepsilon \quad \text{implies} \quad \|\mathcal{N}(x) - \mathcal{N}(x_0)\|_\infty \leq \delta$$

all elements in the input space close to  $x_0$       the classification is roughly the same

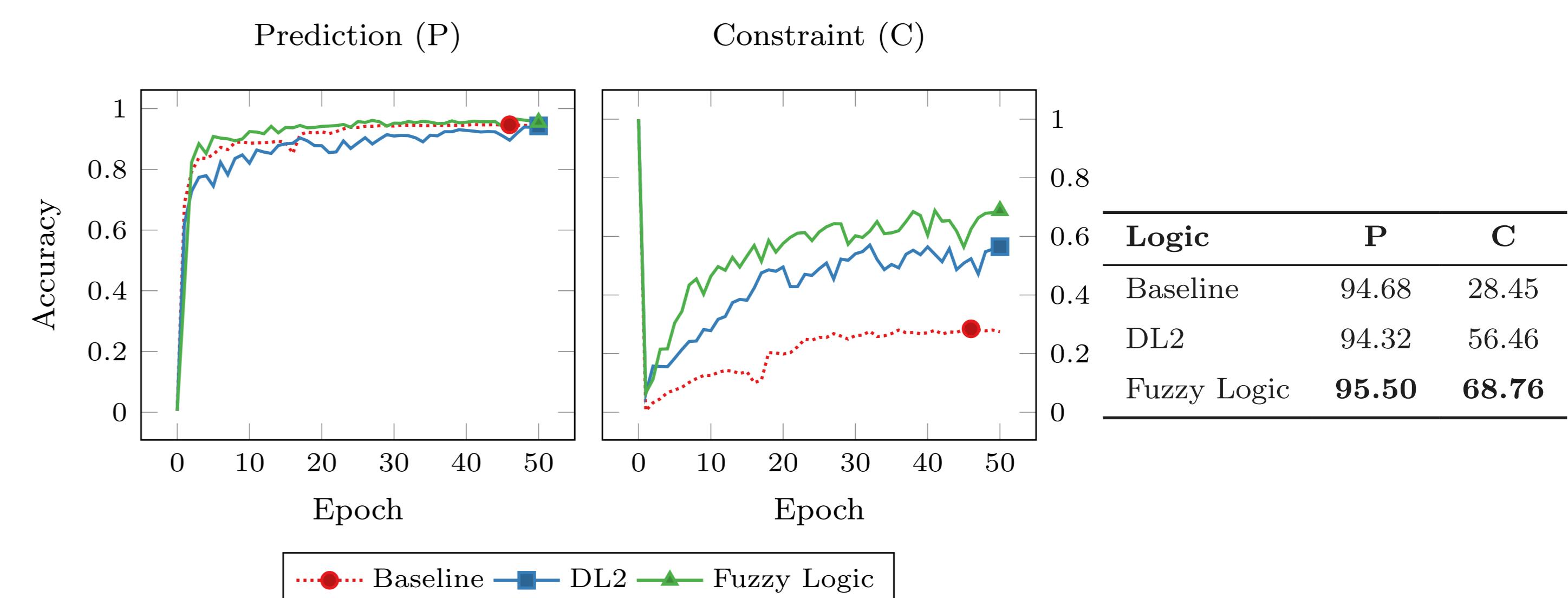
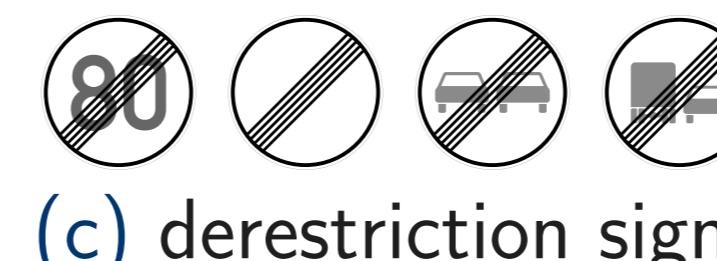


Figure: The Robustness ( $\epsilon = 0.4, \delta = 0.01$ ) constraint on GTSRB.

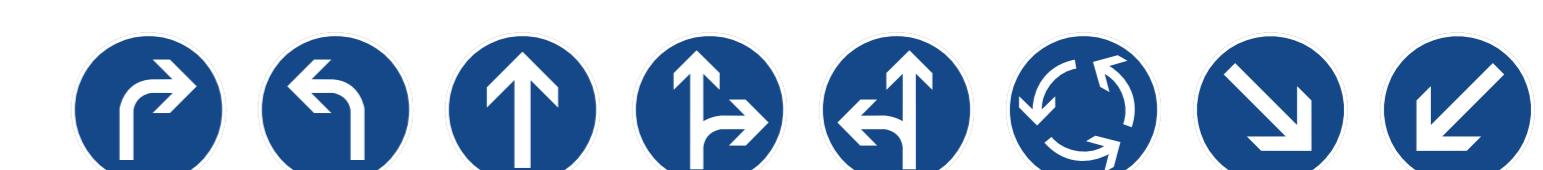
## 6. Experiment – Group Constraint



(b) danger signs



(d) speed limit signs



(f) mandatory signs

**Constraint:** The sum of probabilities of groups of related signs must be either very high or very low.

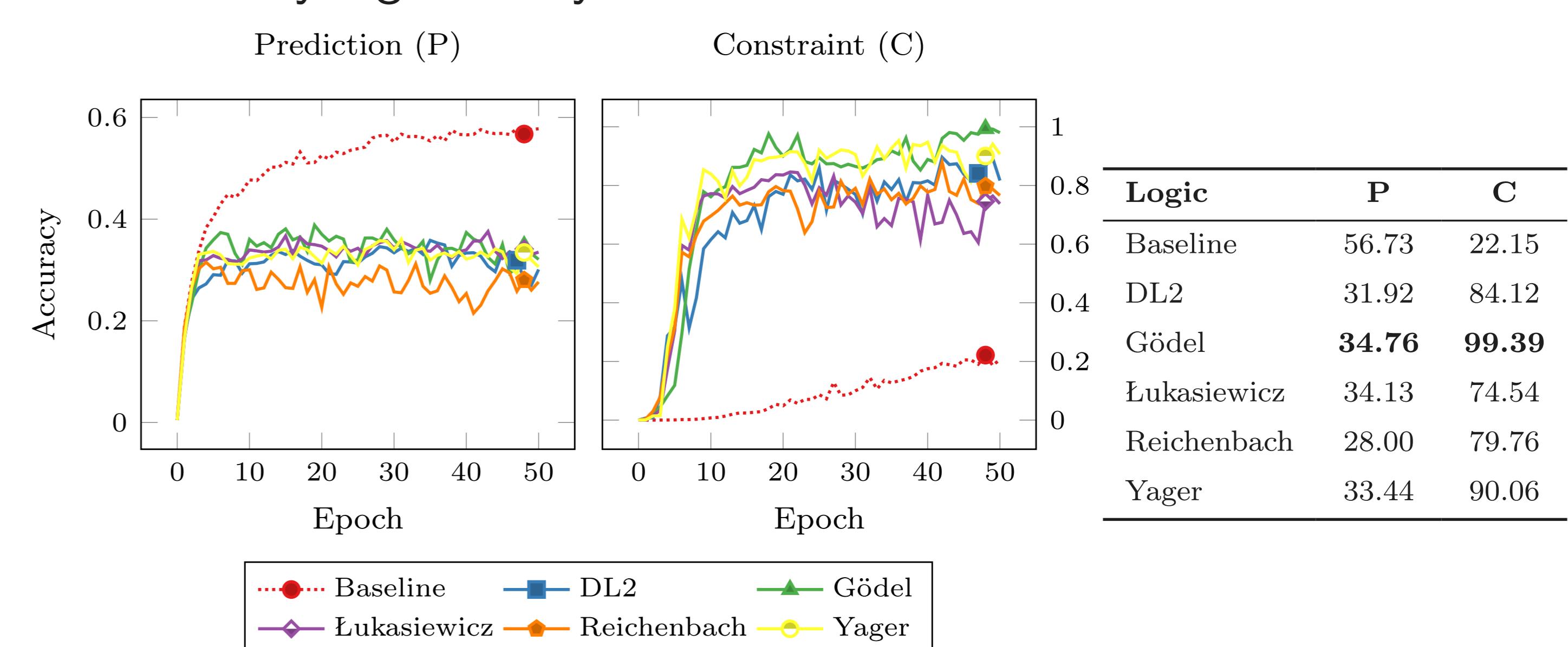


Figure: The Groups ( $\epsilon = 0.6, \delta = 0.02$ ) constraint on GTSRB.

## 7. Results

**Main finding:** Training with *any* differentiable logic works well! Operators with favourable theoretical properties (e.g. shadow-lifting conjunctions, implications with Modus Ponens and Modus Tollens reasoning) are not necessarily the best choice in practice.

## 8. Future Work

- Expressive specifications for ML & temporal differentiable logics.
- Certified training for guaranteed constraint satisfaction.

