

# Algoritmi za brzo učenje na neprijateljskim primjerima

Autor: Dominik Jambrović

Mentor: prof. dr. sc. Siniša Šegvić

# Sadržaj

1. Uvod
2. Robusno učenje
3. Eksperimenti
4. Budući rad
5. Diskusija

# Uvod

- sigurnost umjetne inteligencije
- neprijateljski primjeri
- zatrovani podatci

# Robusno učenje

## Klasični načini

- učenje metodom FGSM
- učenje metodom PGD

## Novi načini

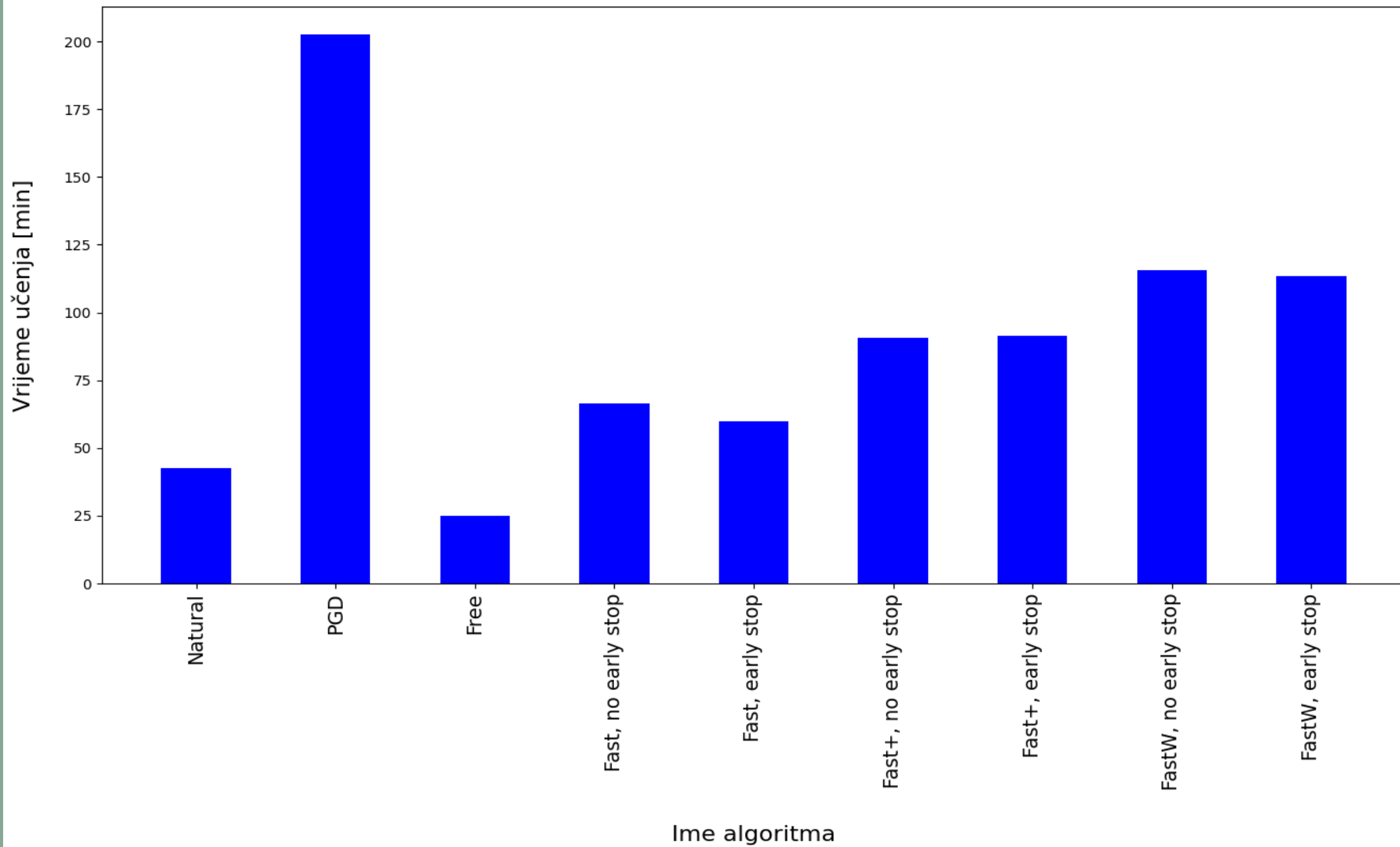
- „besplatno” učenje
- brzo učenje
- nadogradnje brzog učenja

# Eksperimenti

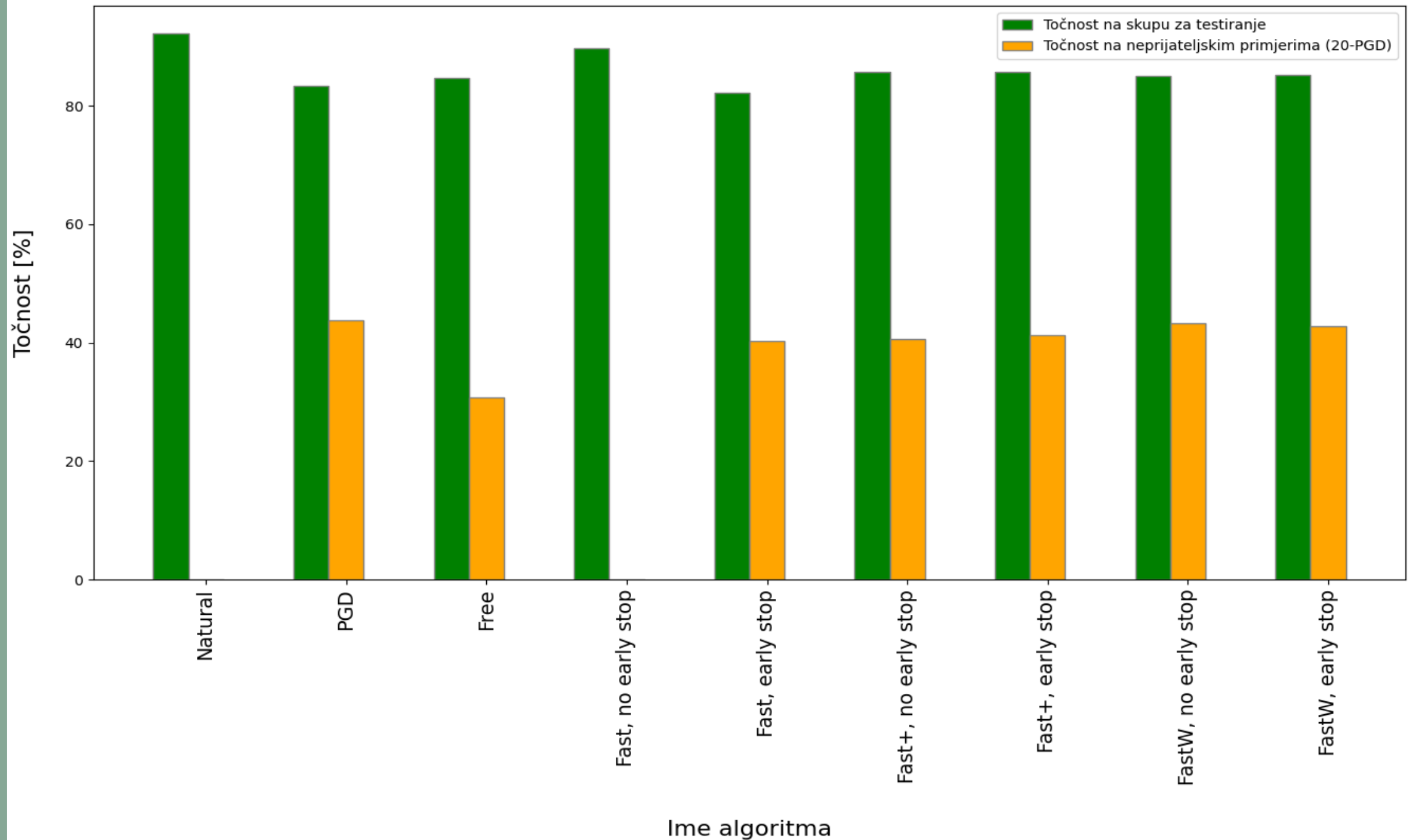
# Izvođenje eksperimenata

- skup podataka CIFAR-10
- arhitektura ResNet-18
- korišćenje računanja u mješovitoj preciznosti
- izvođenje eksperimenata na platformi Kaggle (2x NVIDIA T4)

# Vrijeme učenja za različite algoritme



# Usporedba točnosti za različite algoritme





Pravi razred: cat



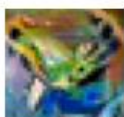
Predvideni razred: dog



Predvideni razred: dog



Predvideni razred: frog



Predvideni razred: dog



Predvideni razred: dog



Predvideni razred: dog



Pravi razred: ship



Predvideni razred: automobile



Predvideni razred: airplane



Predvideni razred: automobile



Predvideni razred: airplane



Predvideni razred: automobile



Predvideni razred: automobile



**Prirodne slike**

Pravi razred: airplane



**Prirodno učenje**

Predvideni razred: ship



**Algoritam PGD**

Predvideni razred: bird



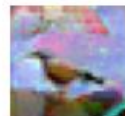
**Algoritam FreeAdv**

Predvideni razred: ship



**Algoritam FastAdv, Early**

Predvideni razred: bird



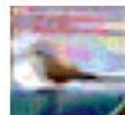
**Algoritam FastAdv+, Early**

Predvideni razred: bird



**Algoritam FastAdvW, Early**

Predvideni razred: bird



Pravi razred: frog



Predvideni razred: automobile



Predvideni razred: deer



Predvideni razred: deer



Predvideni razred: deer



Predvideni razred: cat



Predvideni razred: deer



Pravi razred: automobile



Predvideni razred: truck



Predvideni razred: dog



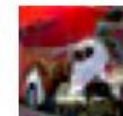
Predvideni razred: truck



Predvideni razred: truck



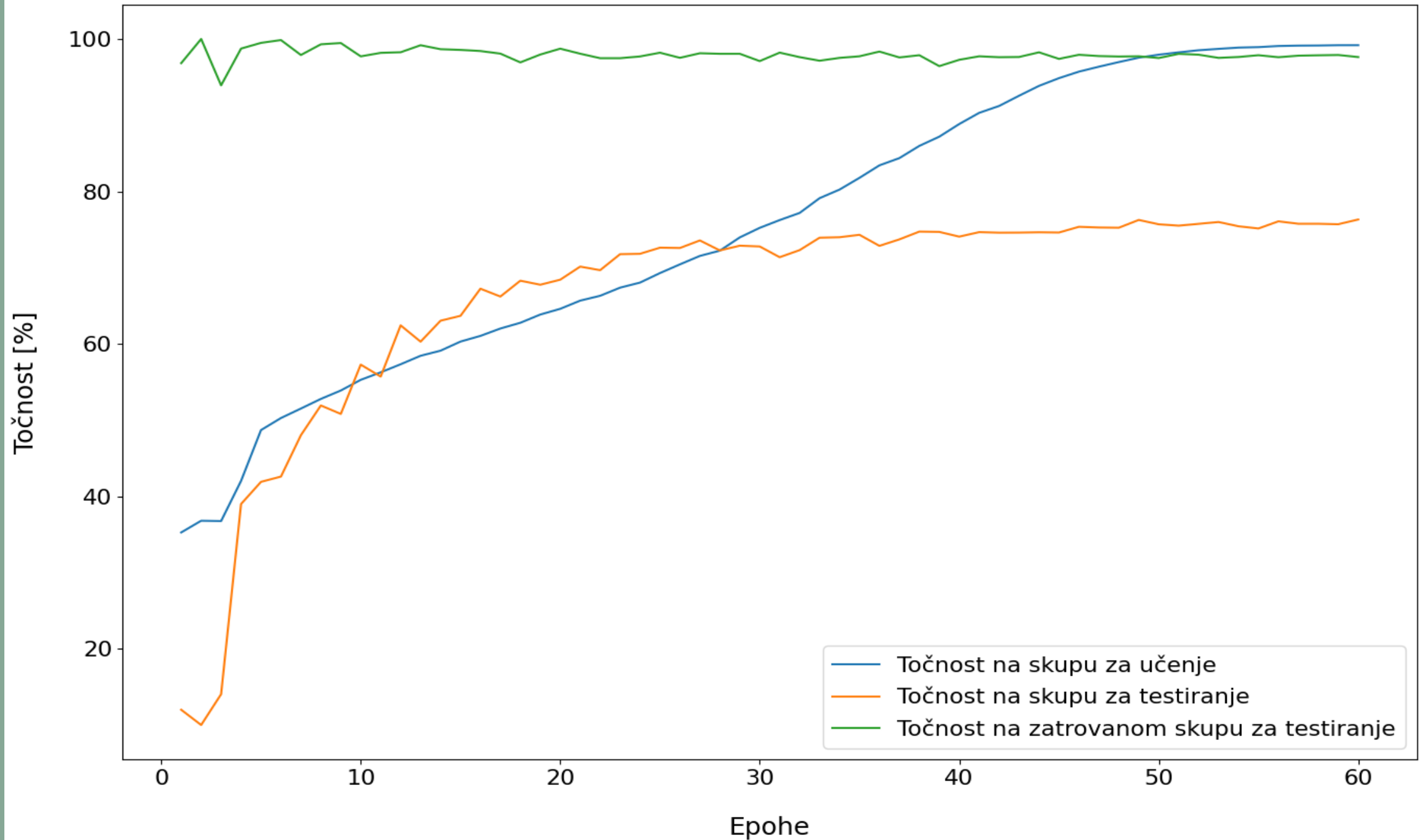
Predvideni razred: truck



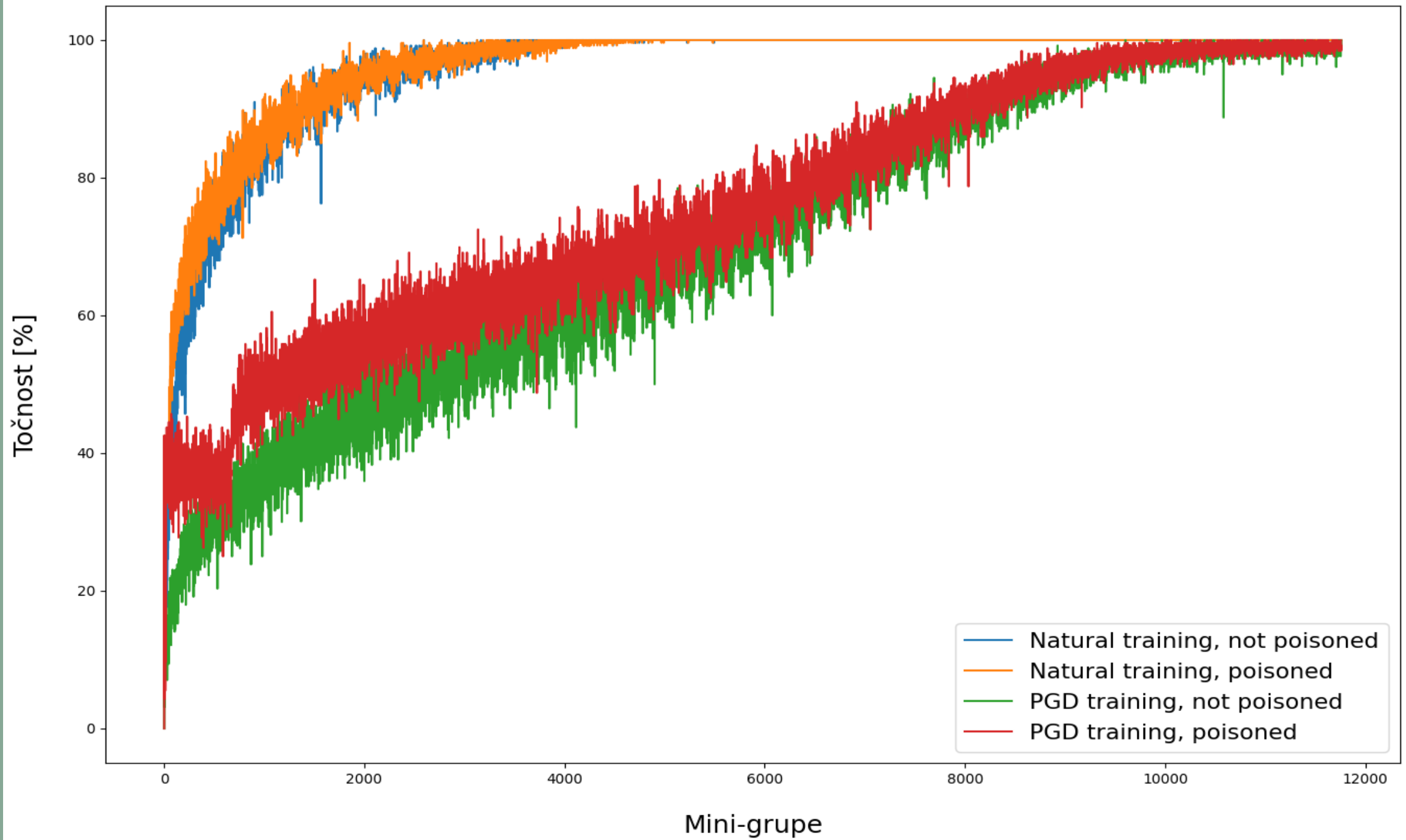
Predvideni razred: cat



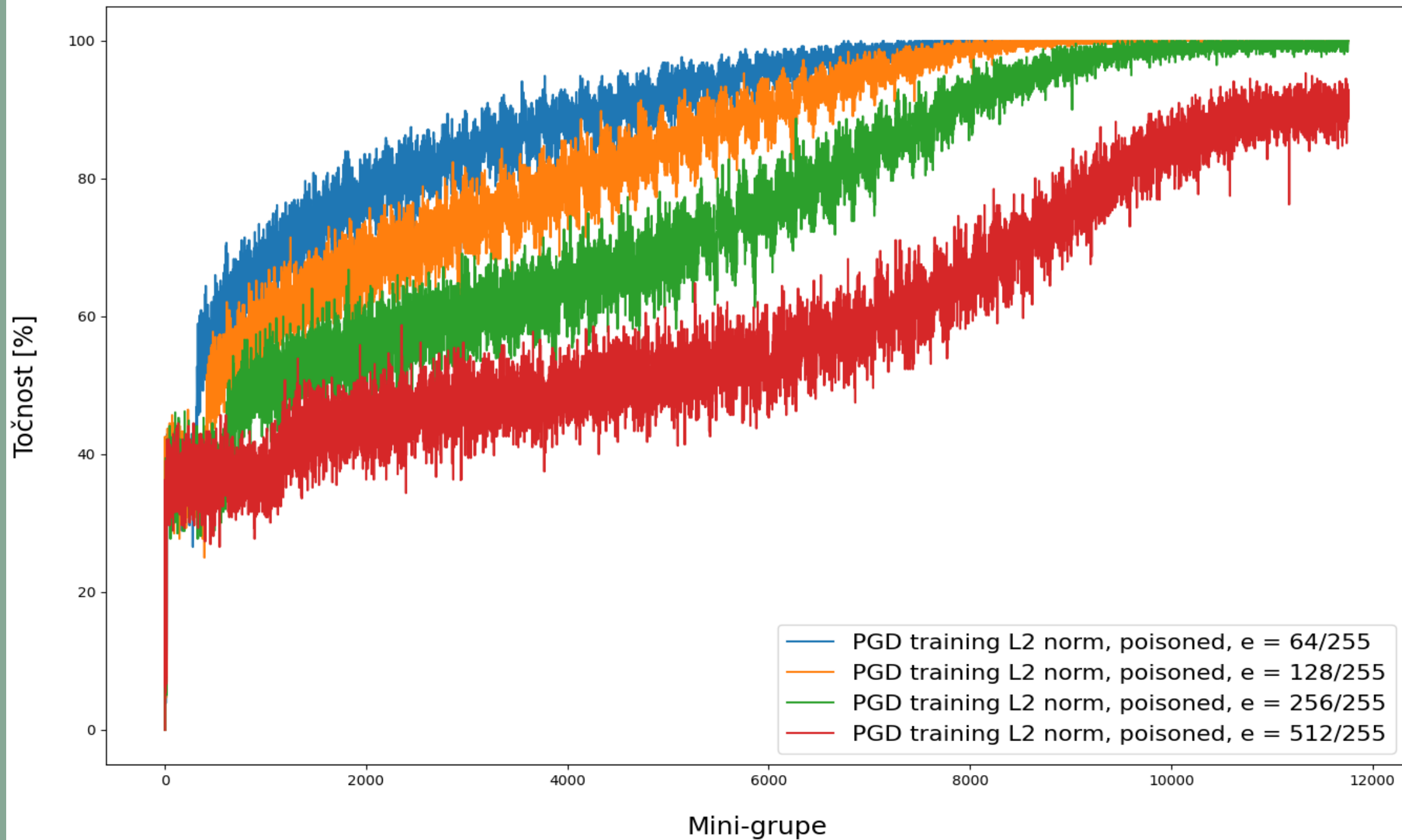
# Usporedba performansi modela Resnet18 PGD, Poisoned



# Točnost po mini-grupama za različite modele

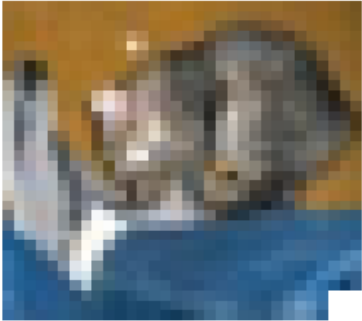


# Točnost po mini-grupama za različite modele



## Zatrovane slike

Pravi razred: cat



Pravi razred: ship

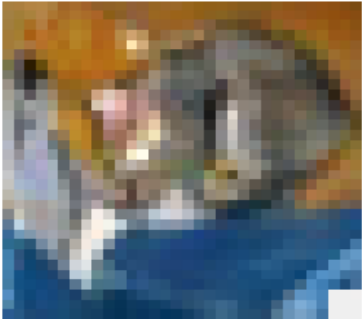


Pravi razred: airplane



## Norma $L_{\infty}$

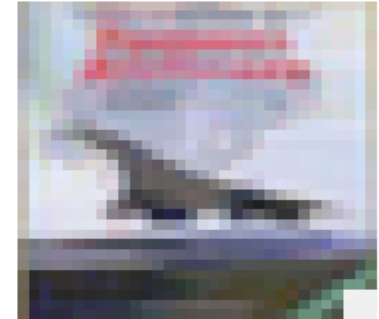
Predviđeni razred: automobile



Predviđeni razred: ship



Predviđeni razred: automobile

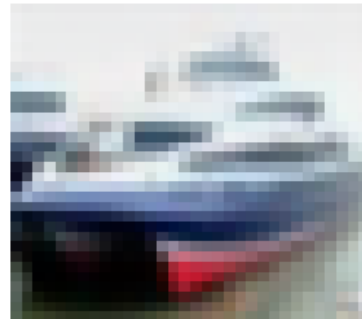


## Norma $L_2$

Predviđeni razred: automobile



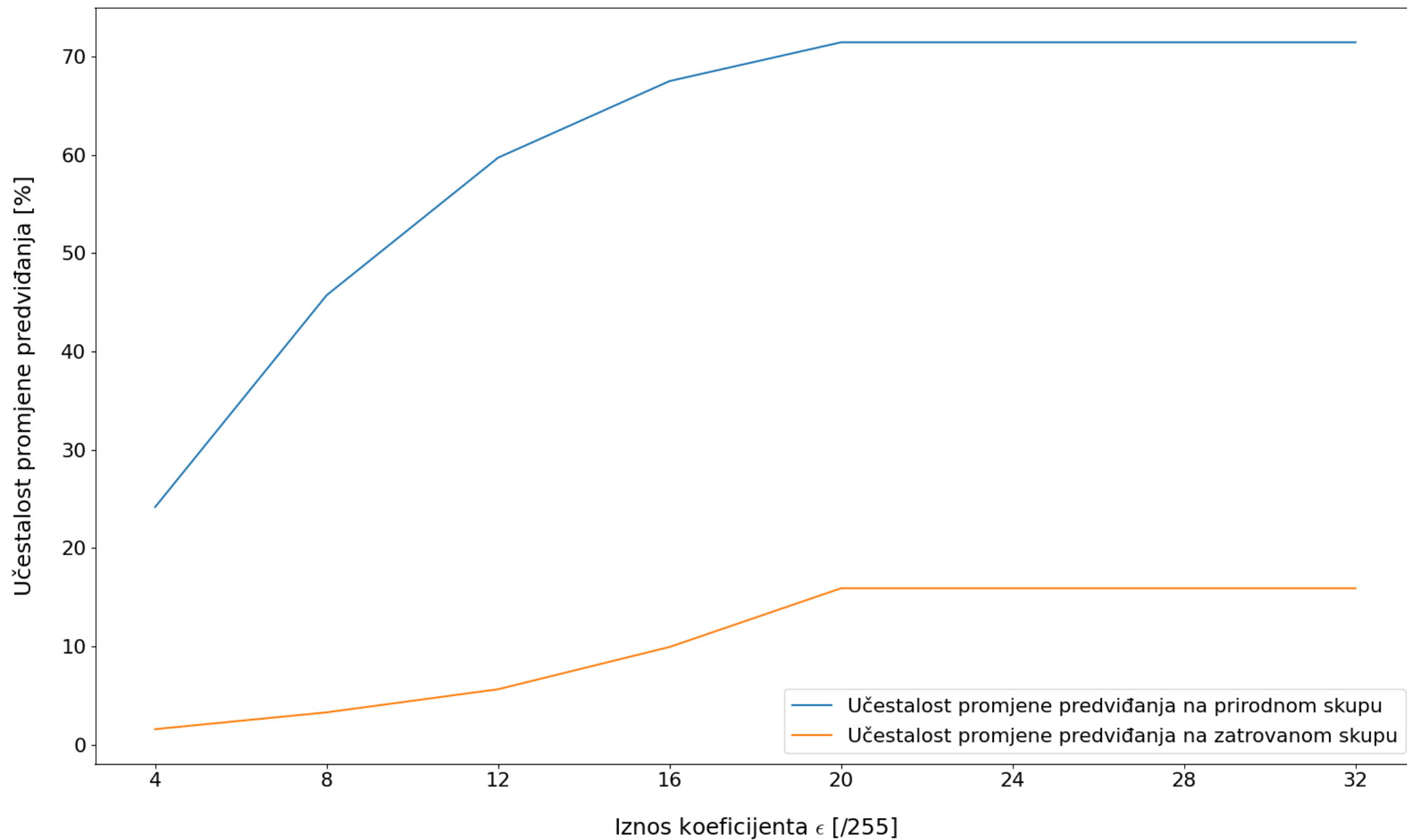
Predviđeni razred: ship



Predviđeni razred: airplane



# Učestalost promjene predviđanja na prirodnom i zatrovanom skupu



# Budući rad

Algoritmi za brzo robusno učenje

- primijeniti metode na kompleksnije arhitekture
- mogućnost kombiniranja „besplatnog” i brzog učenja

Detekcija zatrovanih podataka

- proučiti utjecaj korištenja L1 norme
- mjeriti uspješnost detekcije zatrovanih podataka praćenjem promjena predviđanja

# Diskusija