



# Federated Learning under Attack: Improving Gradient Inversion for Batch of Images



**FURG**  
UNIVERSIDADE FEDERAL  
DO RIO GRANDE

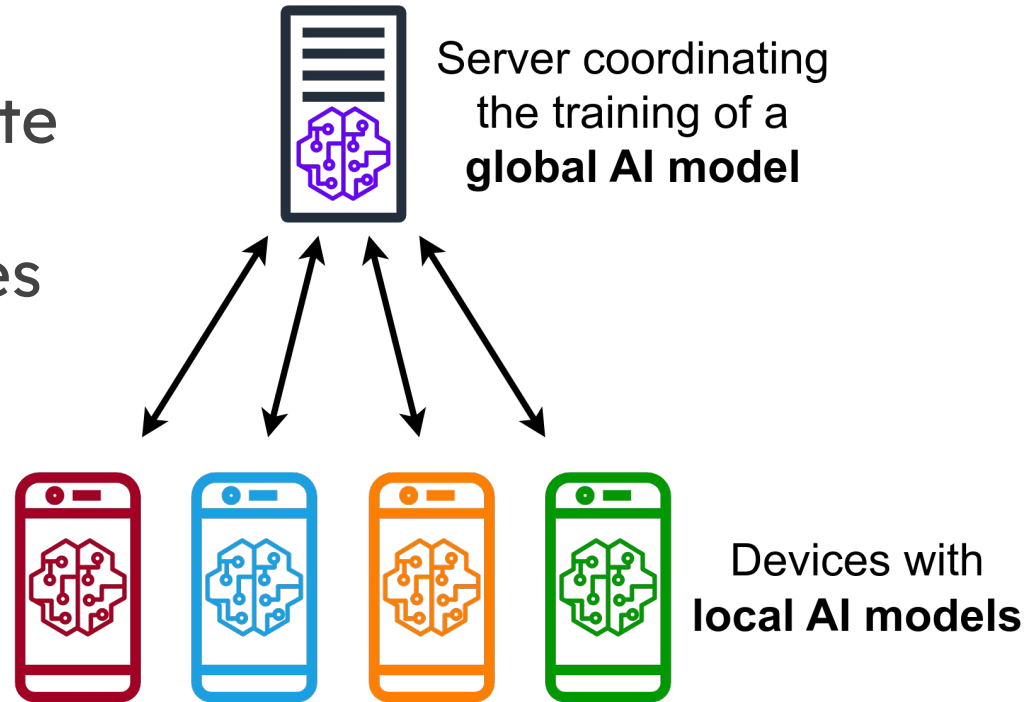


Luiz Leite, Yuri Santos  
Prof. Bruno Dalmazo,  
Orientador: Prof. André Riker

UFPA  
FURG

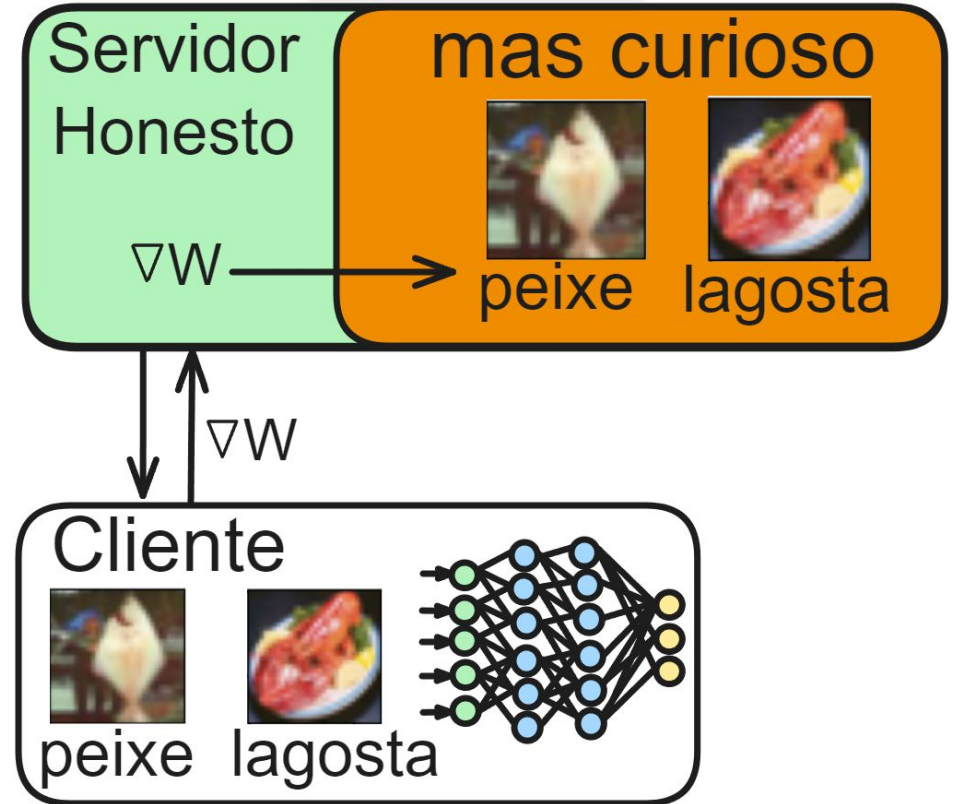
# Motivação

- Tecnologia emergente
- Ameaças emergentes
- Dados sensíveis



# Problema(s)

- Ameaças emergentes:
  - Inversão de gradiente
  - Servidor honesto, mas curioso
  - DLG e iDLG



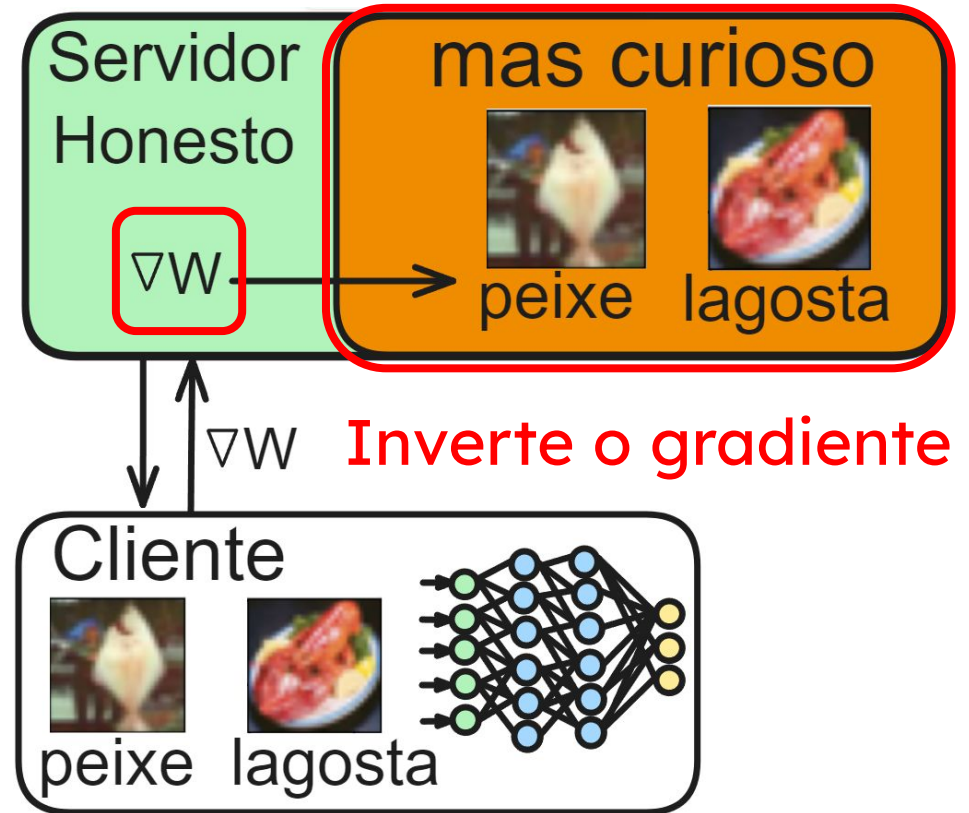
# Problema(s)

- Ameaças emergentes:
  - Inversão de gradiente
  - Servidor honesto, mas curioso
  - DLG e iDLG



# Problema(s)

- Ameaças emergentes:
  - Inversão de gradiente
  - Servidor honesto, mas curioso
  - DLG e iDLG



# Problema(s)

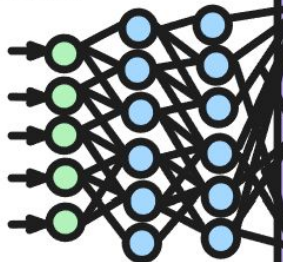
- Ameaças emergentes:
  - Inversão de gradiente
  - Servidor honesto, mas curioso
  - DLG e iDLG



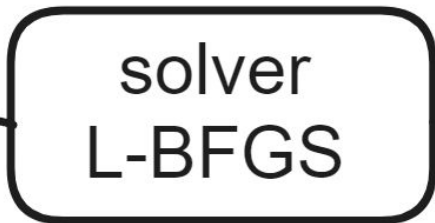
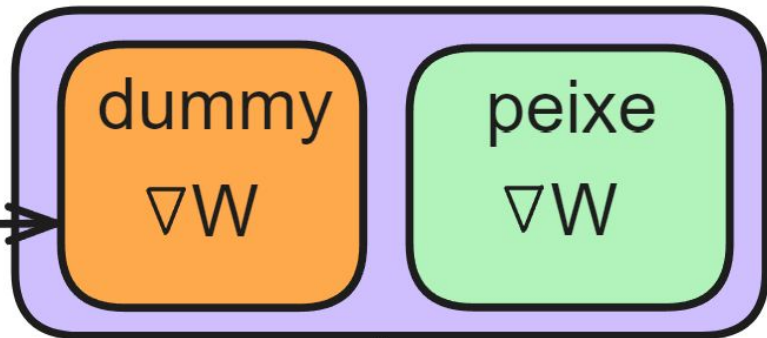
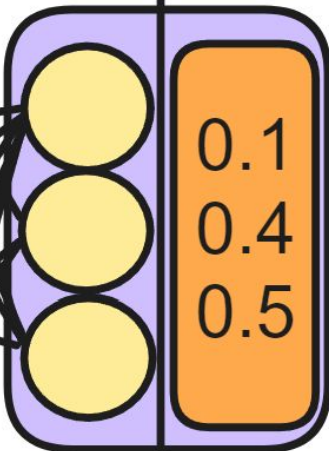
# Servidor

Iteração 1

dummy data

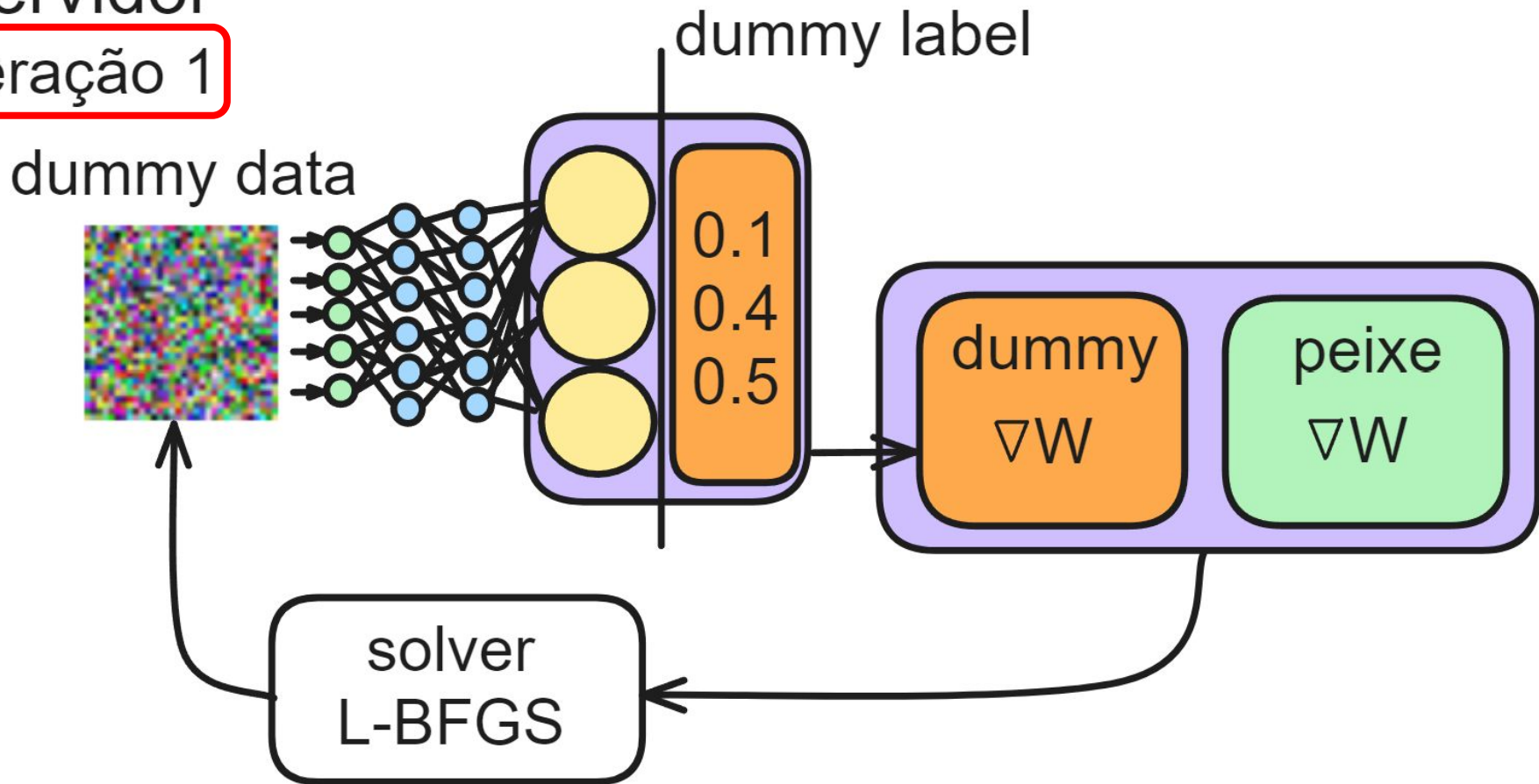


dummy label



# Servidor

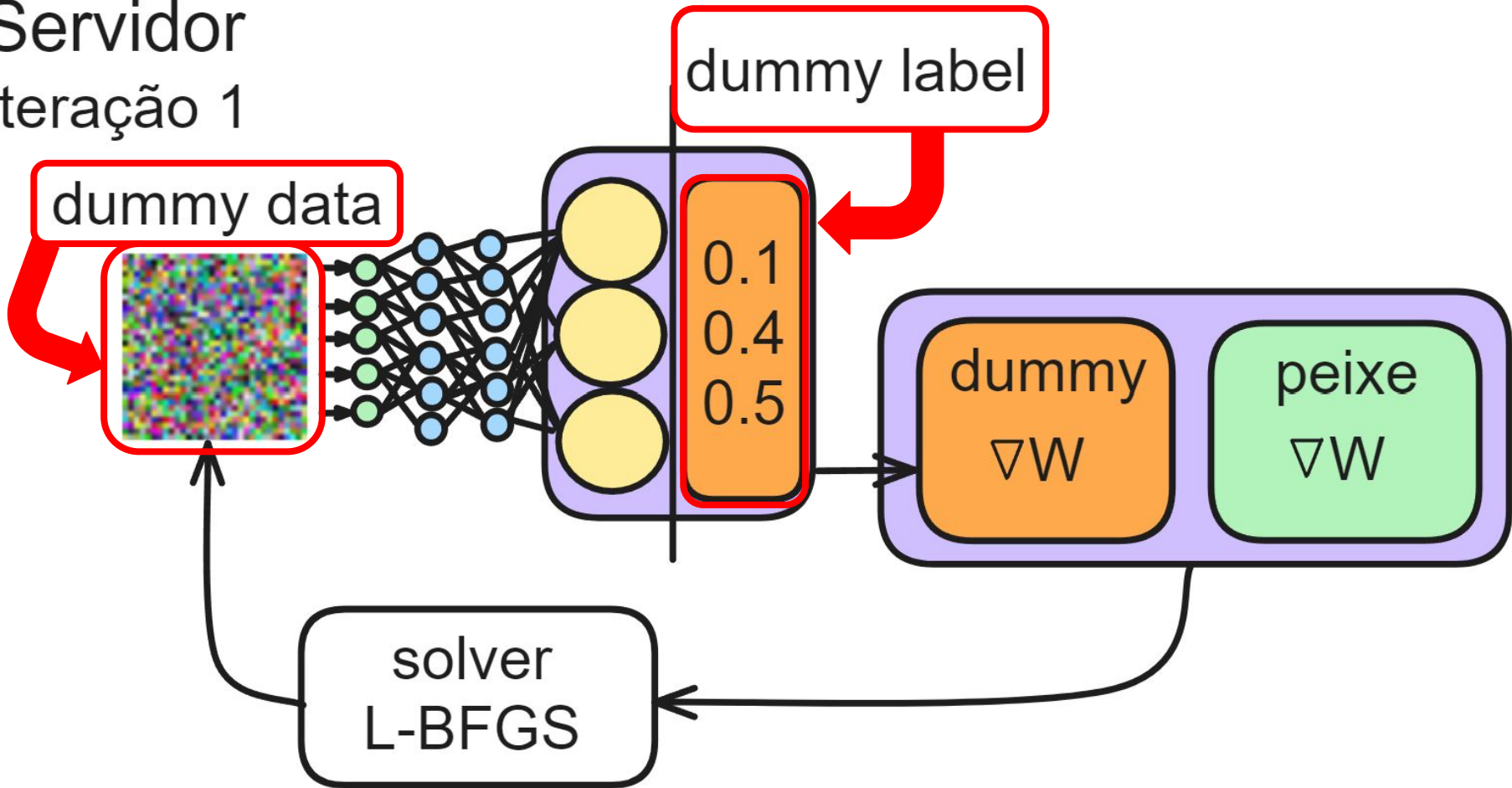
Iteração 1





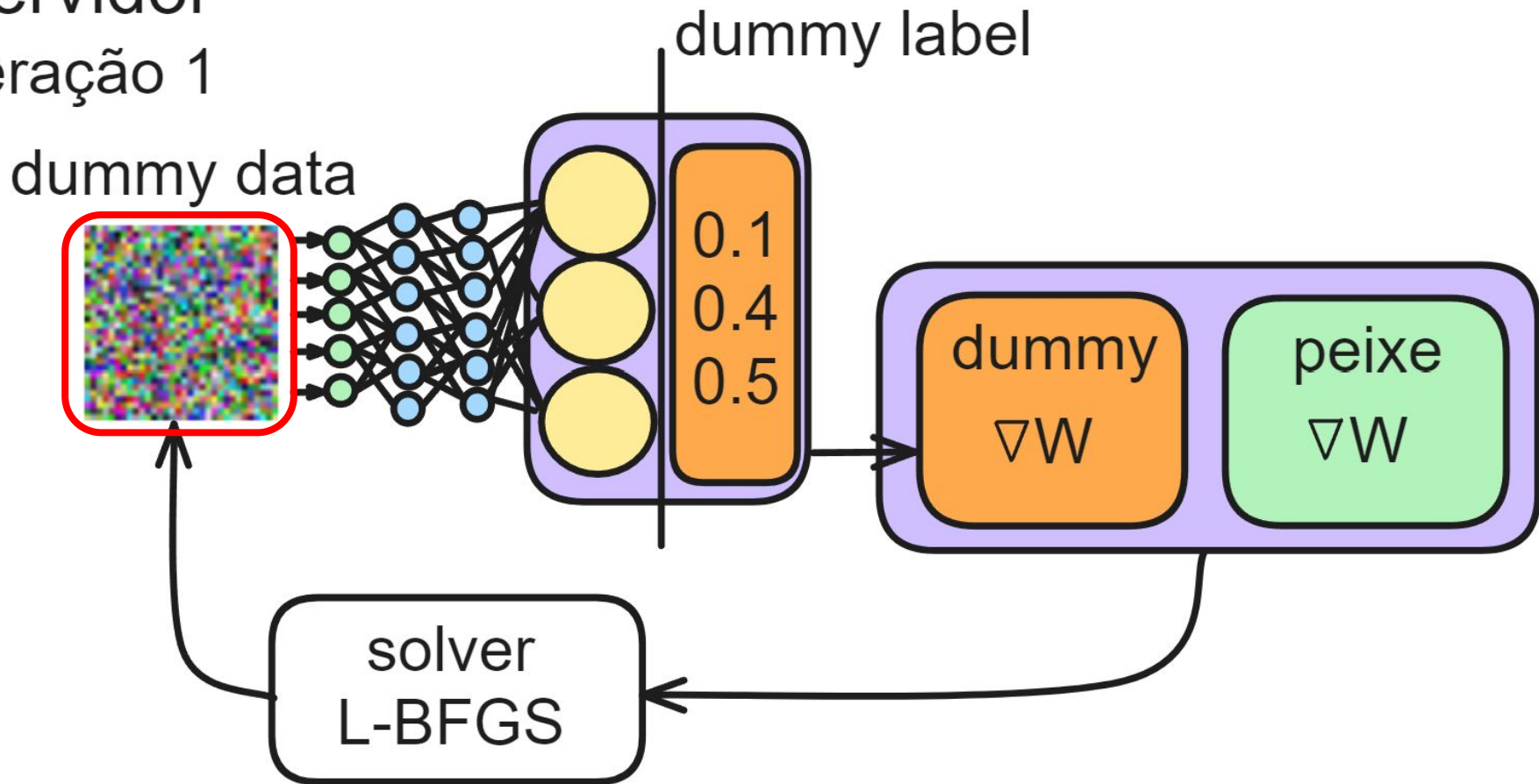
# Servidor

Iteração 1



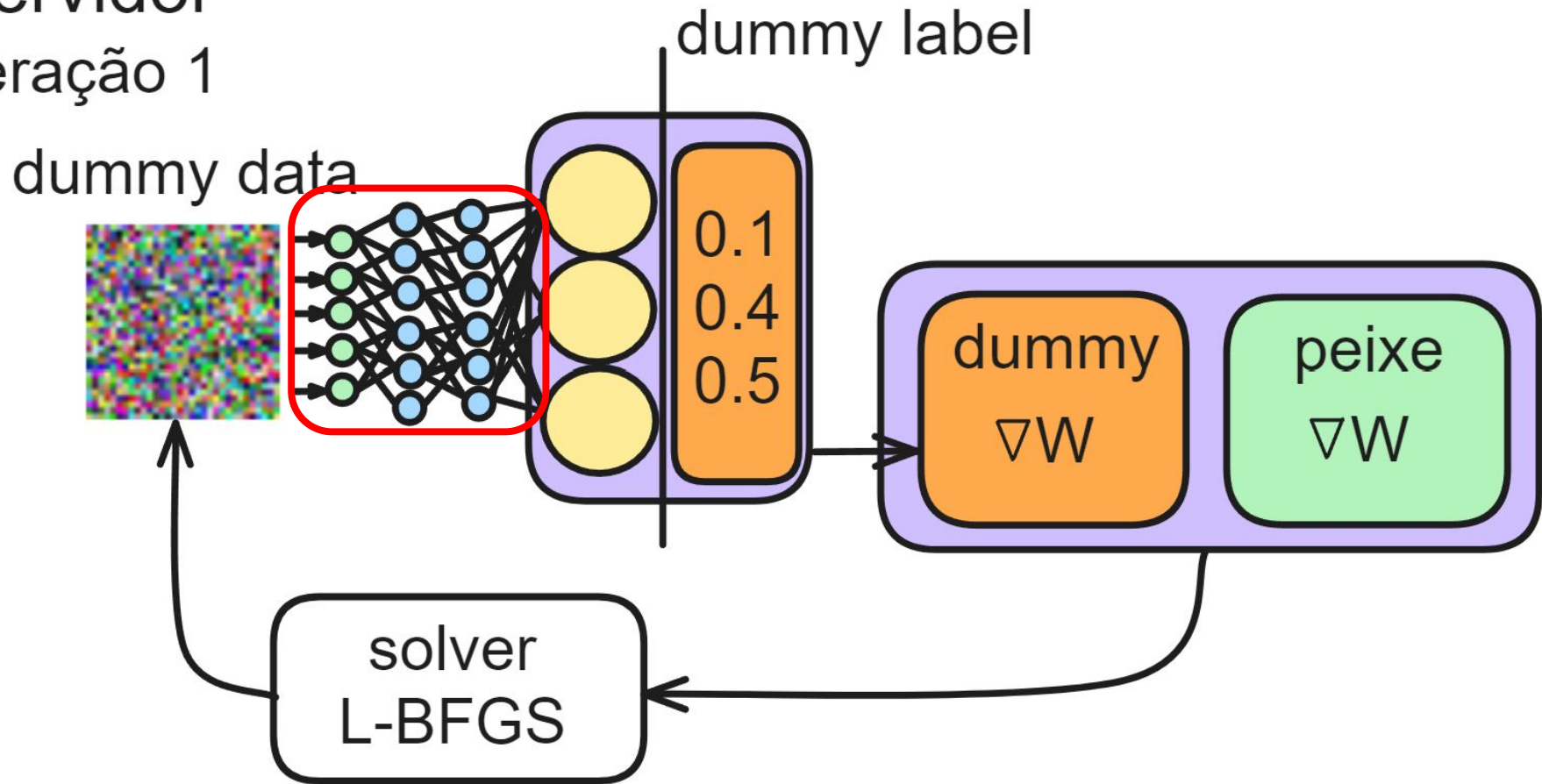
# Servidor

Iteração 1



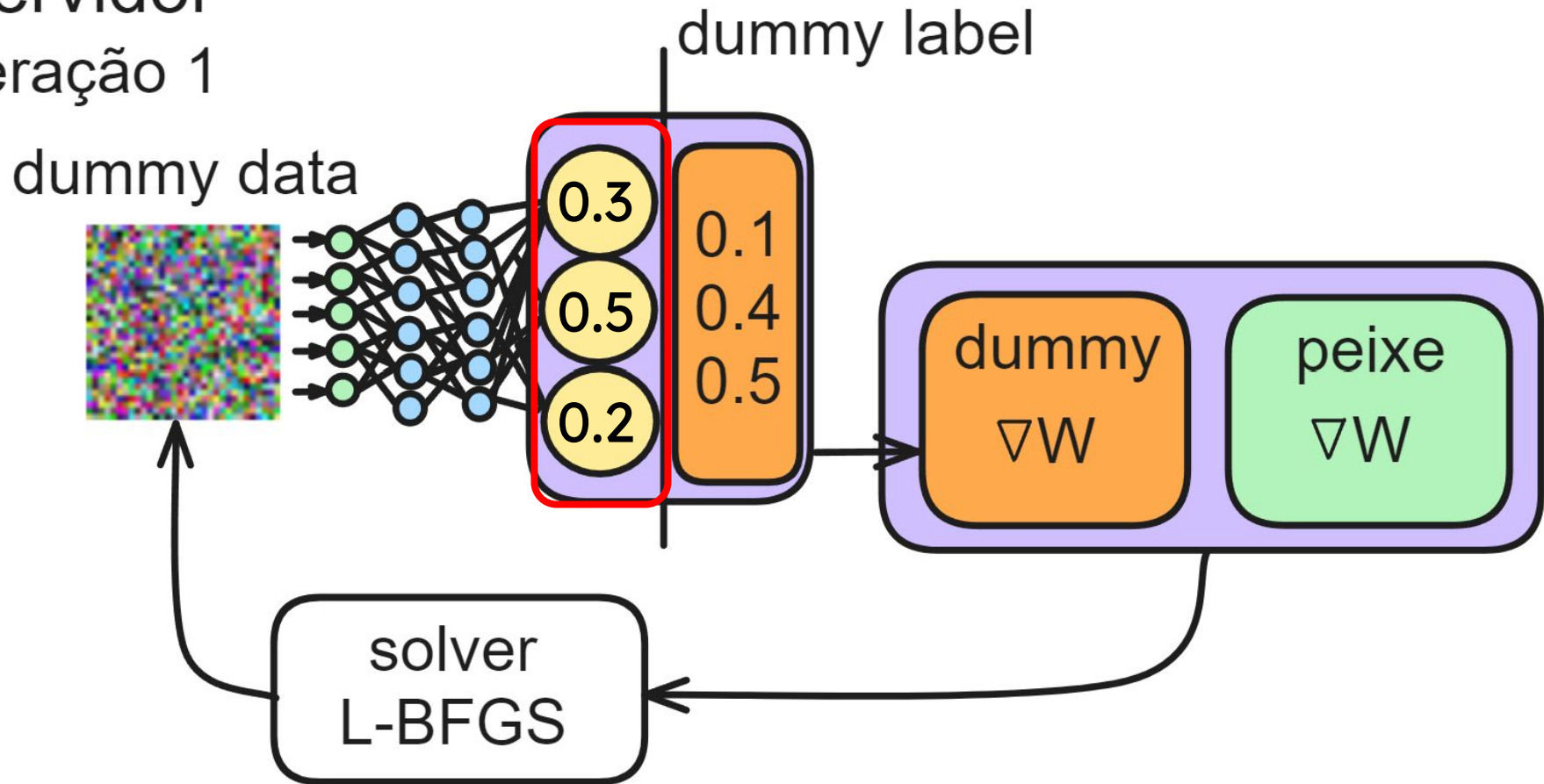
# Servidor

## Iteração 1



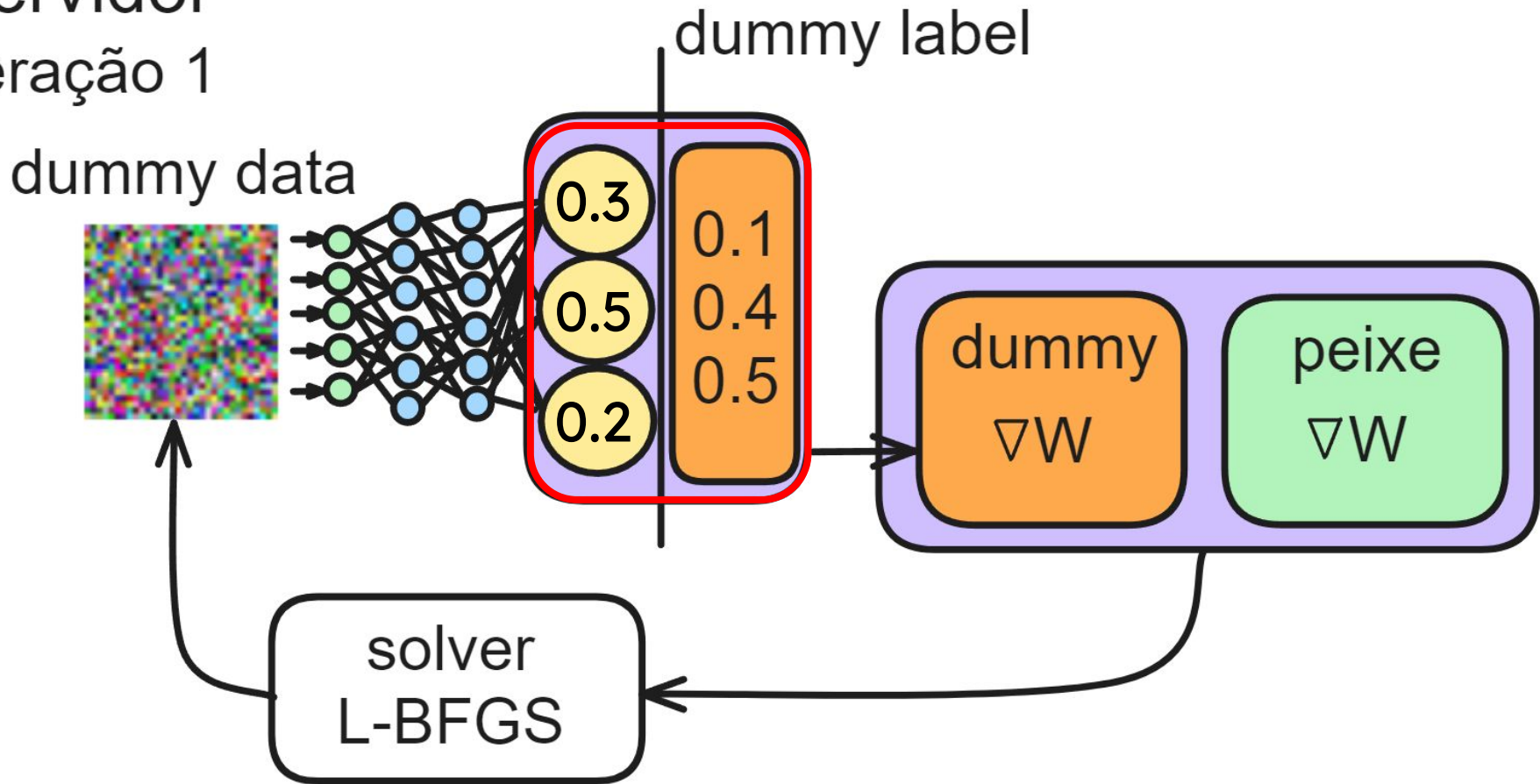
# Servidor

Iteração 1



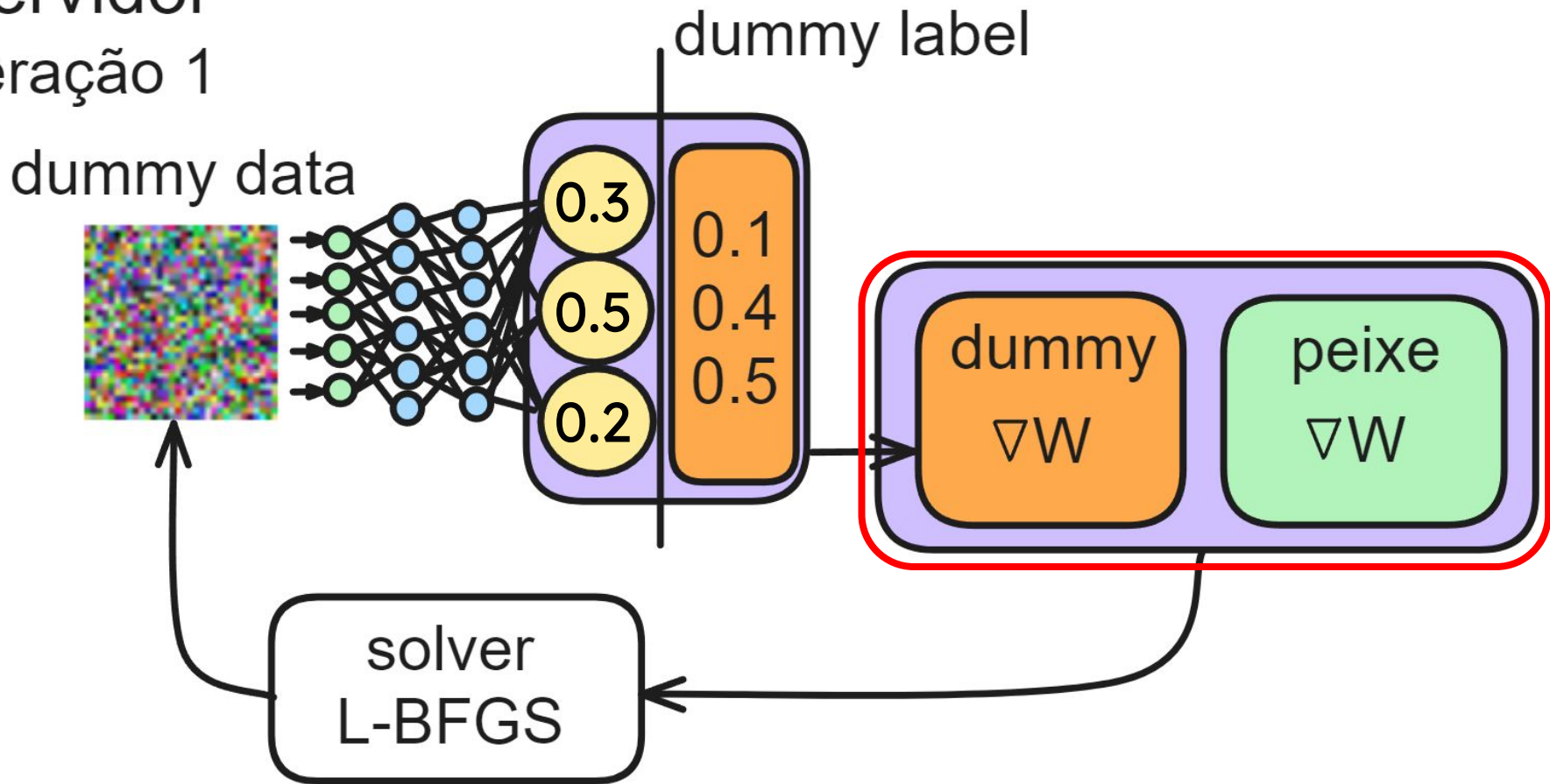
# Servidor

Iteração 1



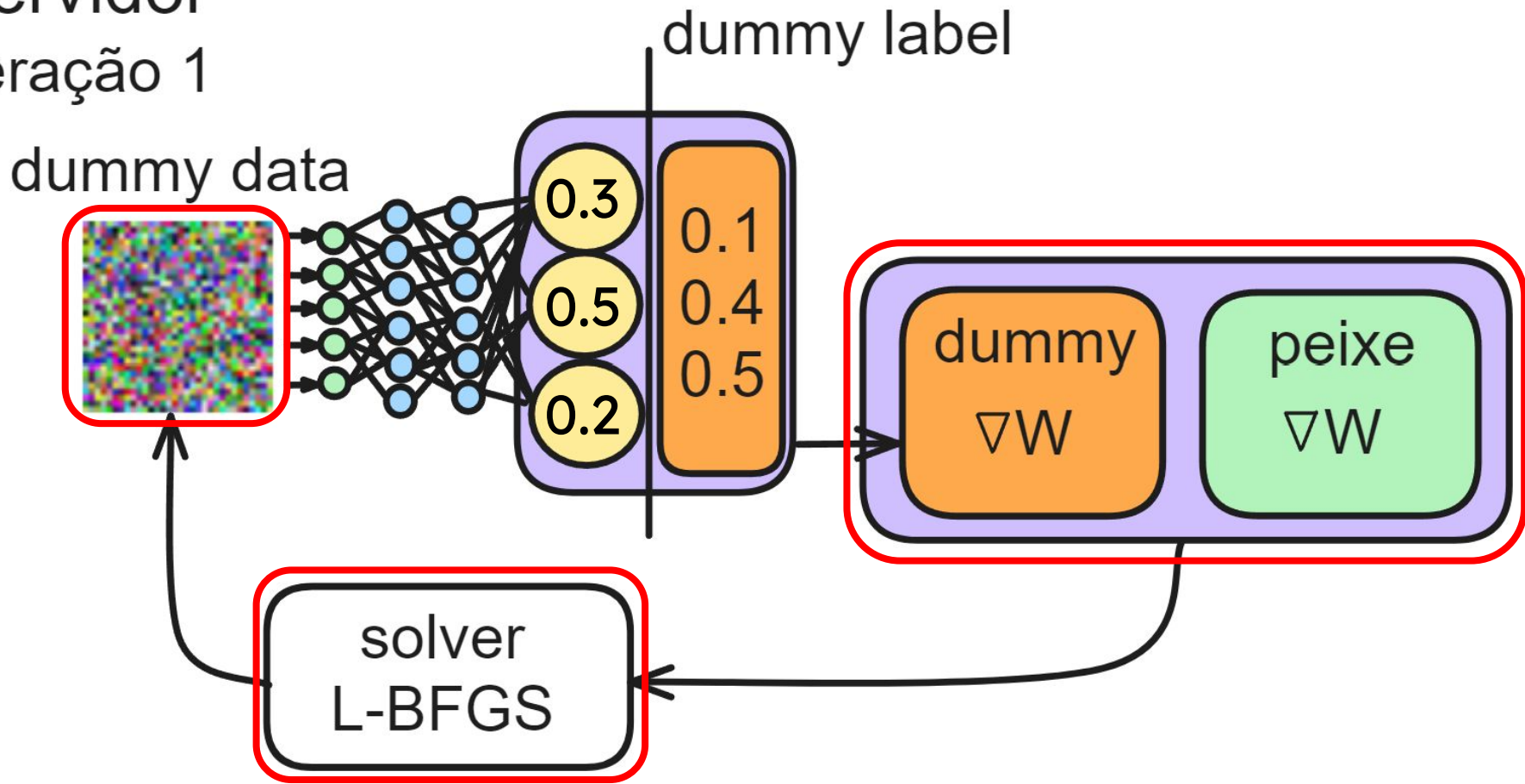
# Servidor

Iteração 1



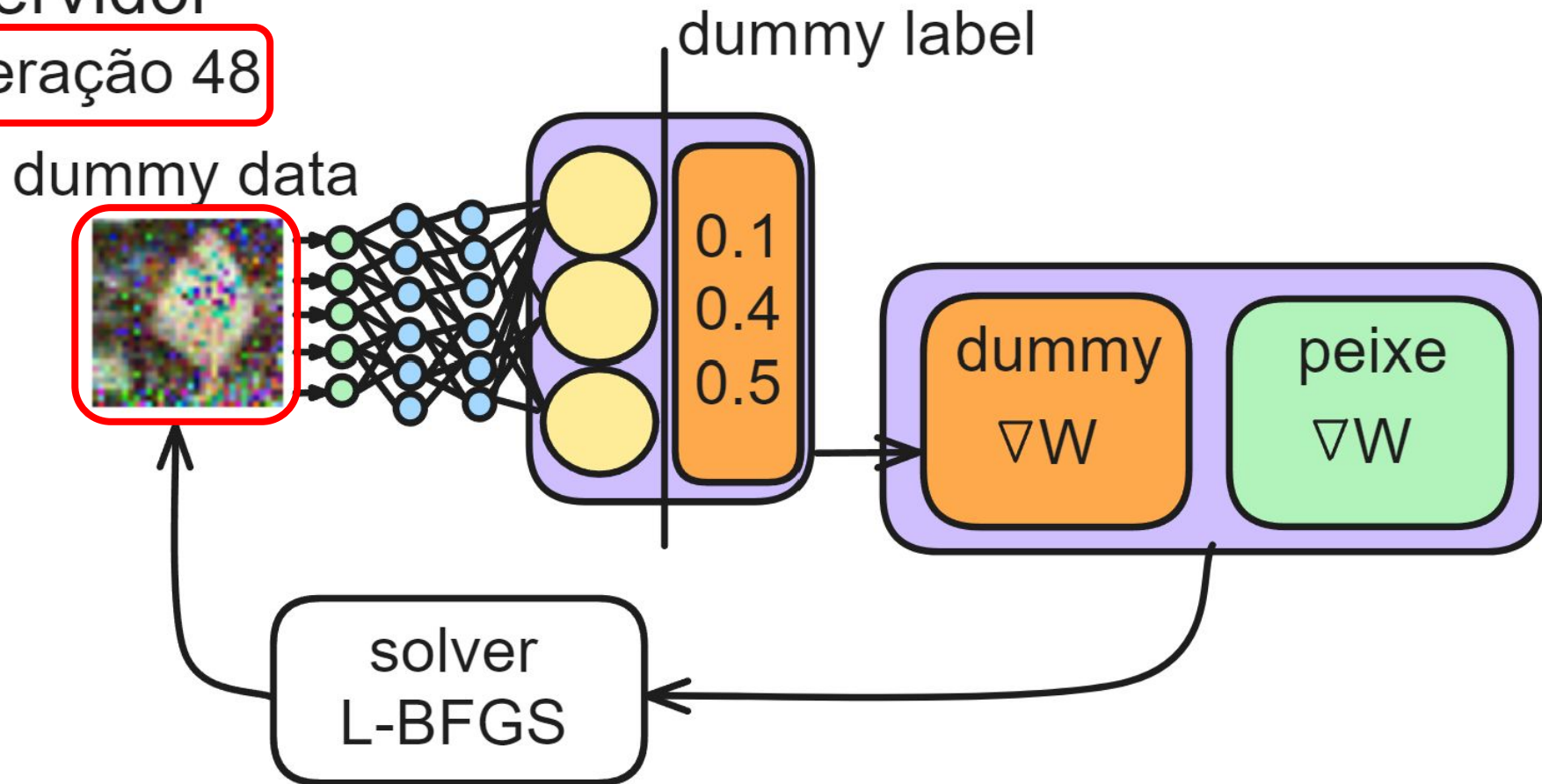
# Servidor

Iteração 1



# Servidor

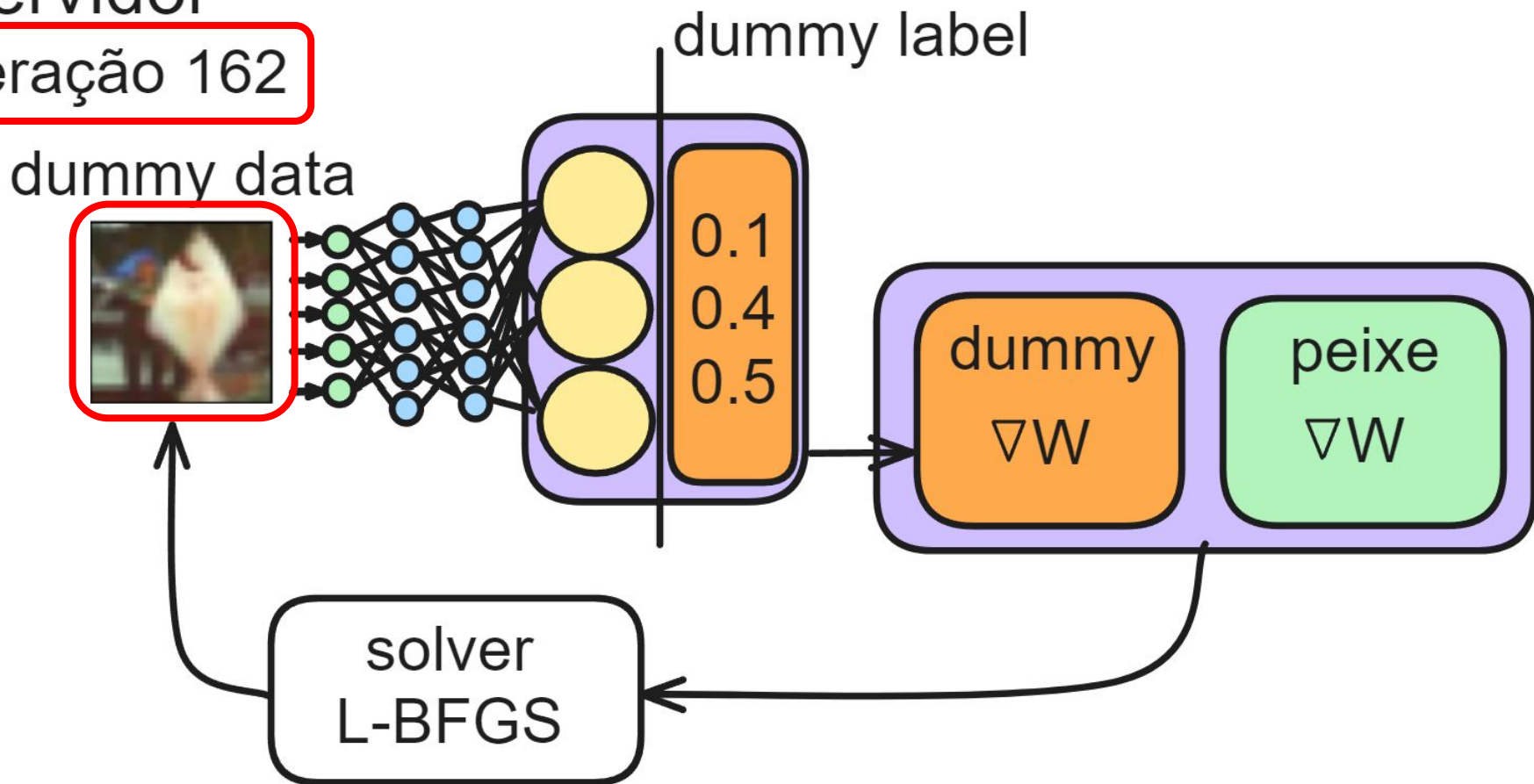
Iteração 48





# Servidor

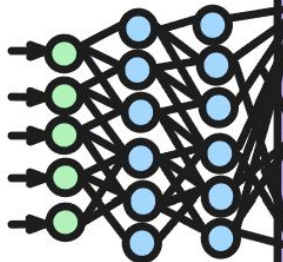
Iteração 162



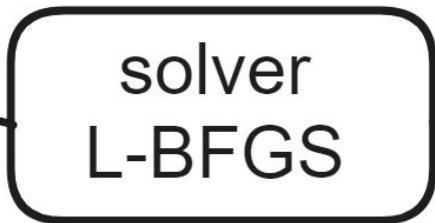
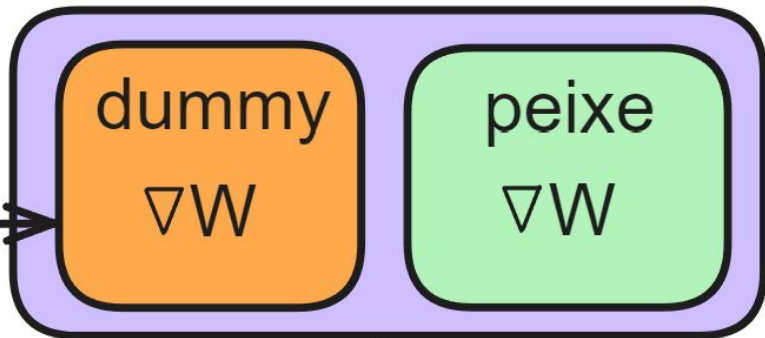
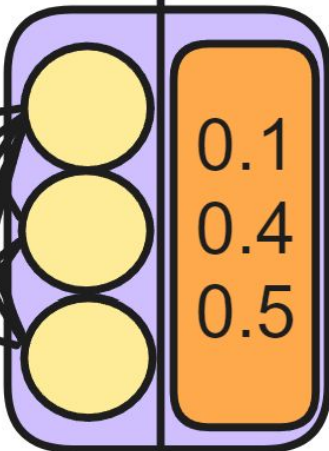
# Servidor

Iteração ???

dummy data



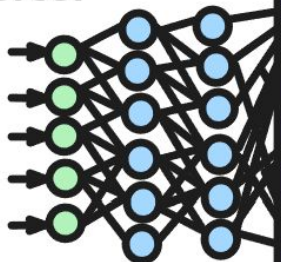
dummy label



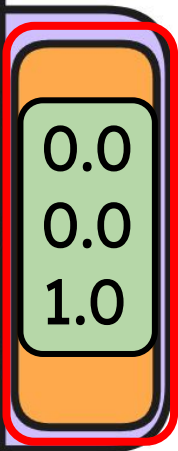
# Servidor

Iteração ???

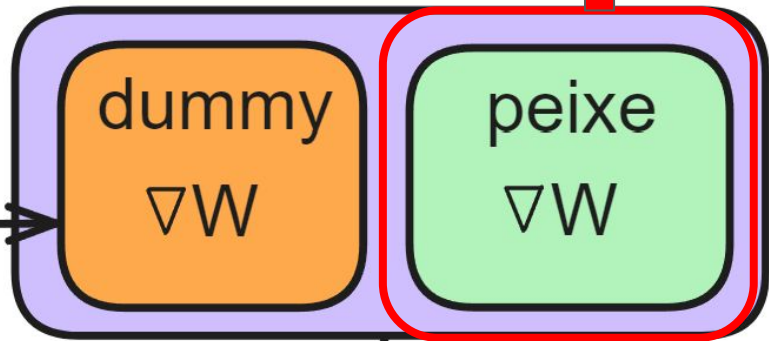
dummy data



dummy label



iDLG



# Desafio(s)

- Defesas
- Quanto perigoso é esse ataque?

# Desafio(s)

- Defesas
- **Quão perigoso é esse ataque?**

Nova técnica descoberta!

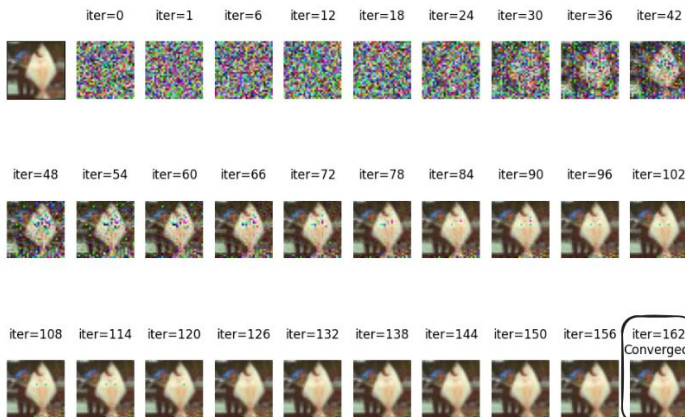


## DLG-FB

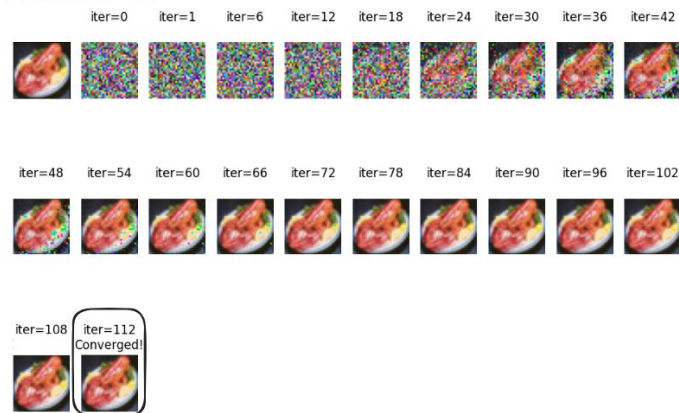


# DLG-FB

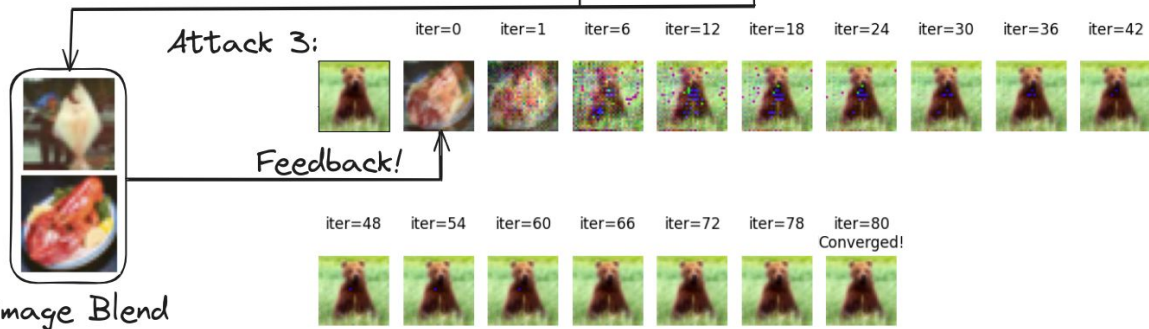
Attack 1:



Attack 2:



Attack 3:



# DLG-FB

Servidor  
Iteração 1

dummy data

dummy label

0.1  
0.4  
0.5

dummy  $\nabla W$

peixe  $\nabla W$

solver  
L-BFGS

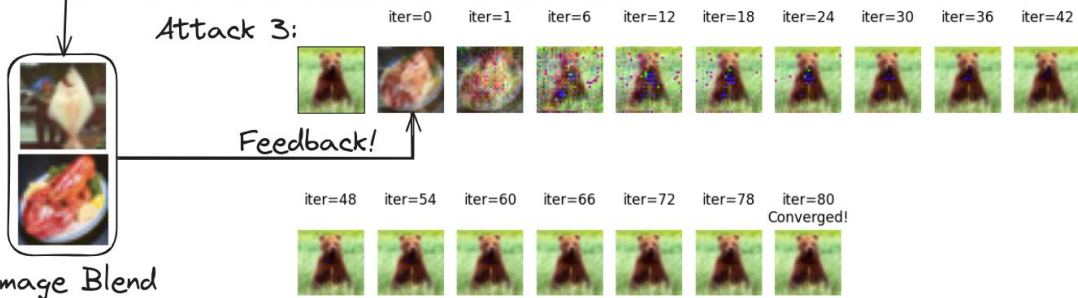
Attack 1:



Attack 2:



Attack 3:



Feedback!

Image Blend

# DLG-FB





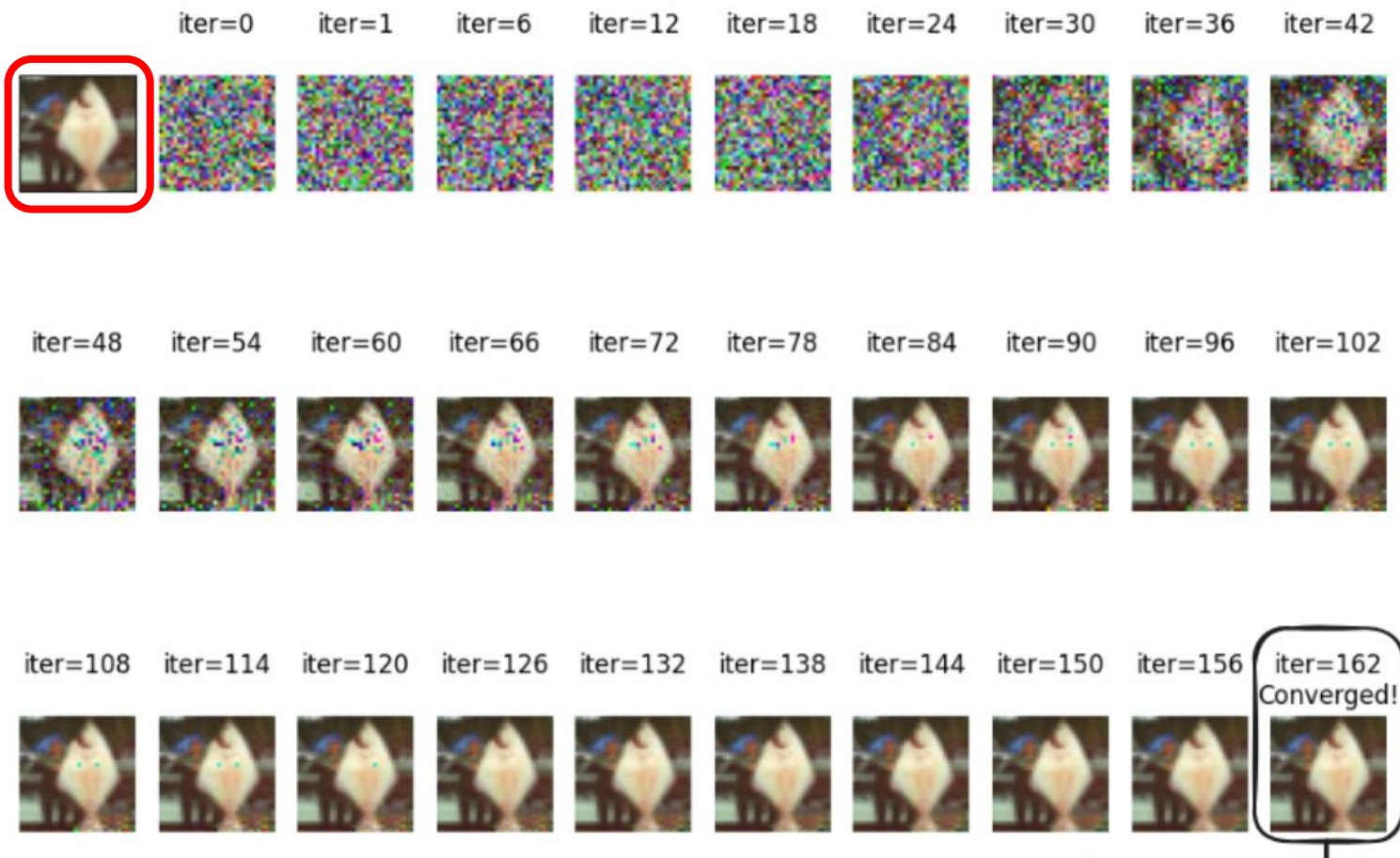
# DLG-FB



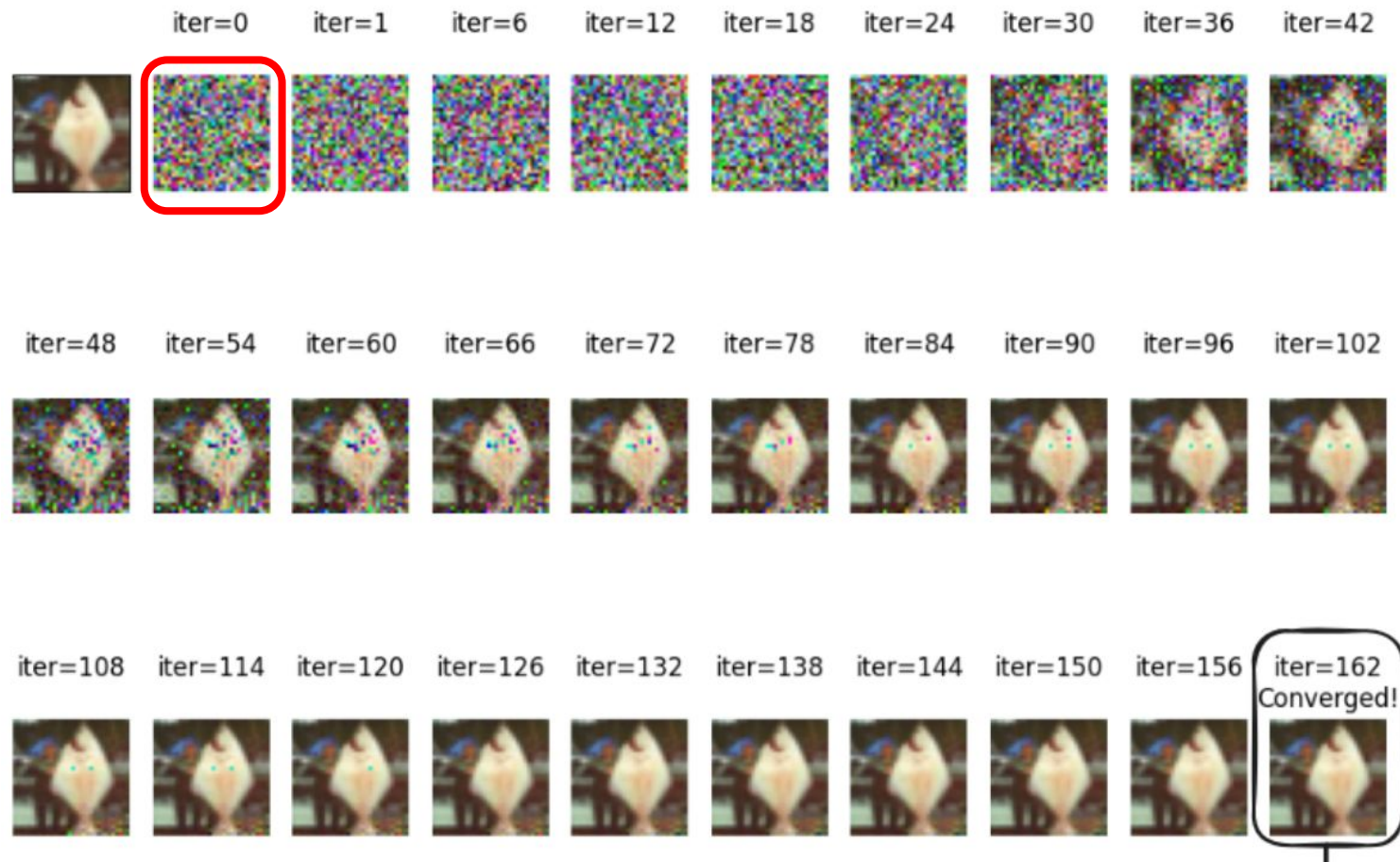
# DLG-FB



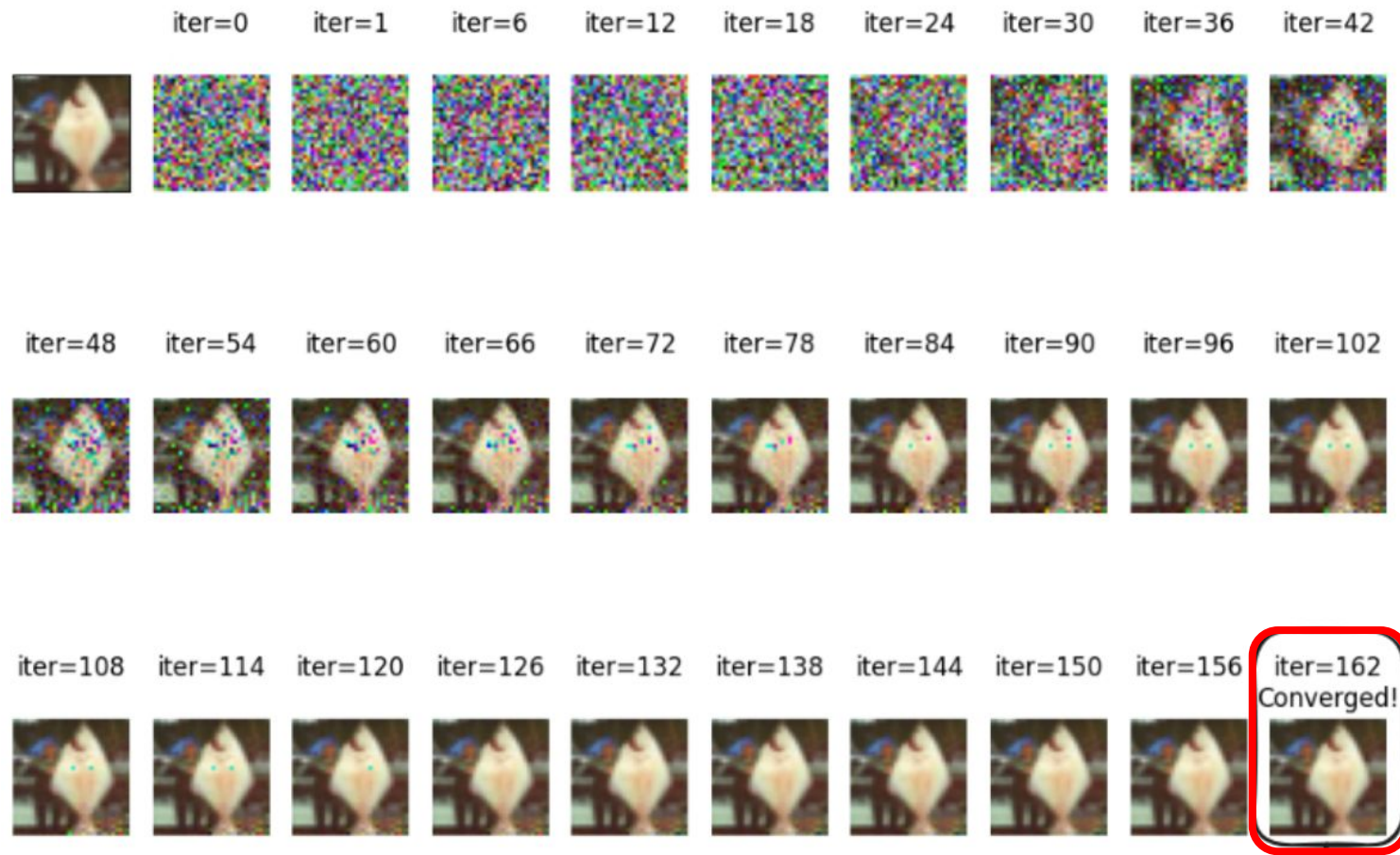
# Attack 1:



# Attack 1:



# Attack 1:



# Attack 2:

iter=0

iter=1

iter=6

iter=12

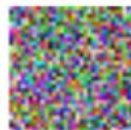
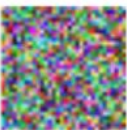
iter=18

iter=24

iter=30

iter=36

iter=42



iter=48

iter=54

iter=60

iter=66

iter=72

iter=78

iter=84

iter=90

iter=96

iter=102



iter=108

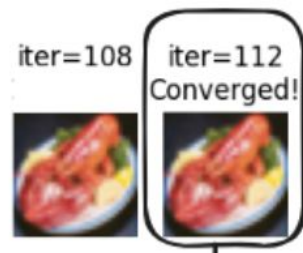
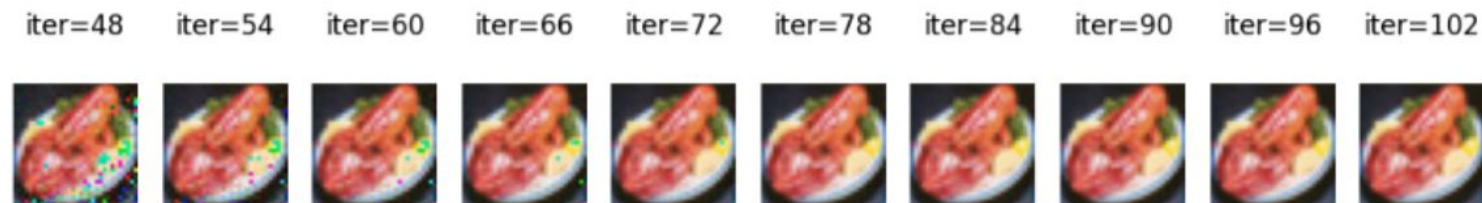
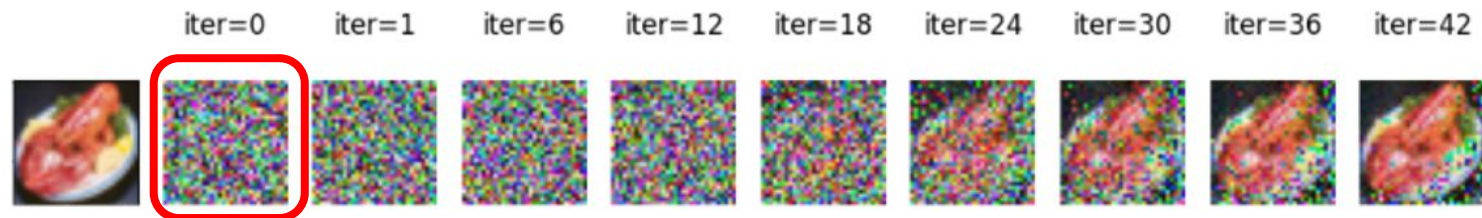


iter=112

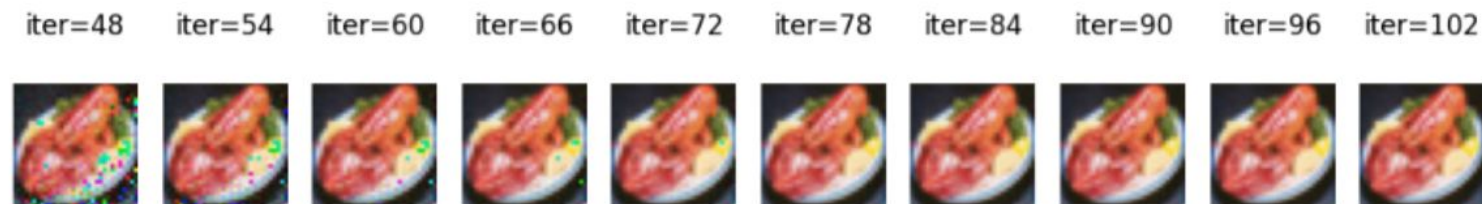
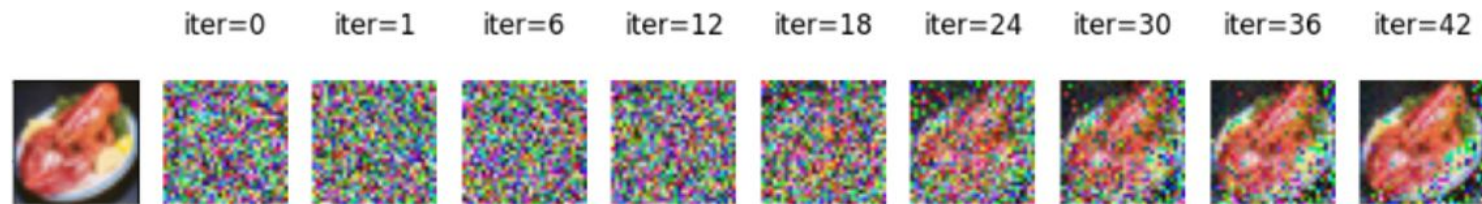
Converged!



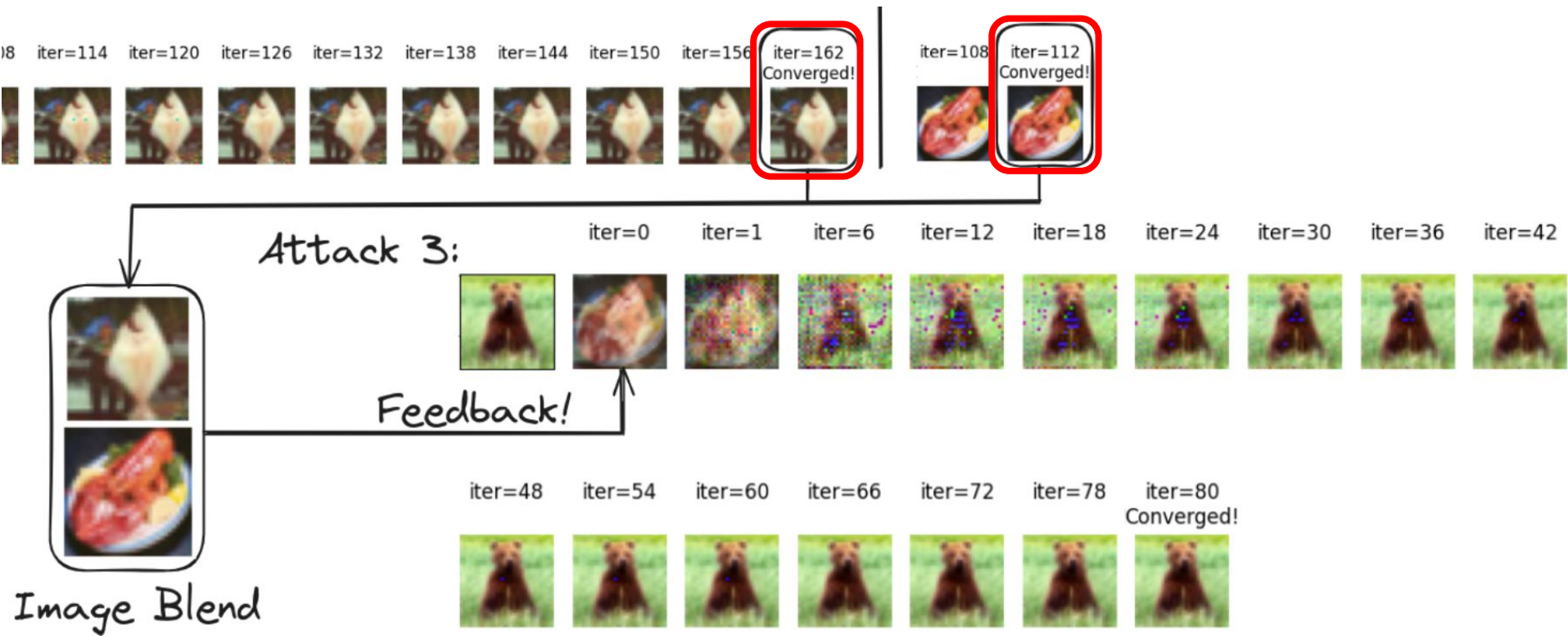
# Attack 2:

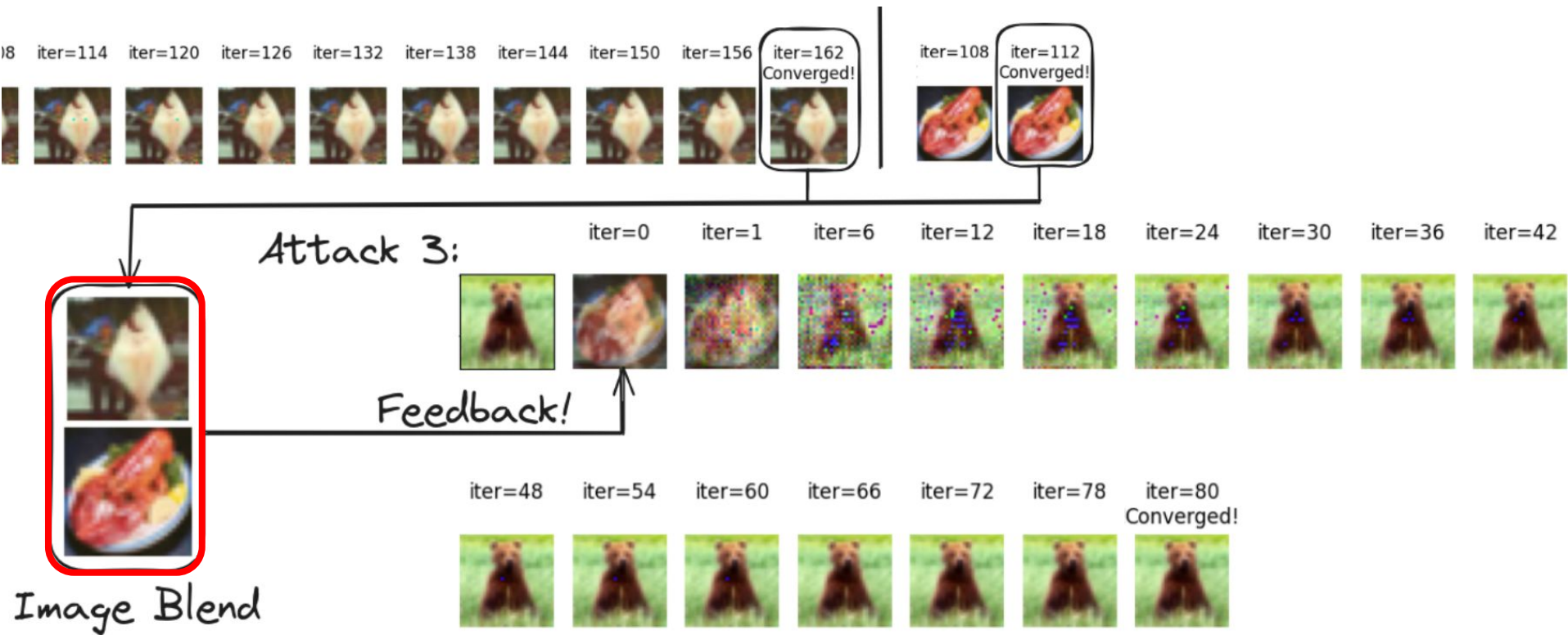


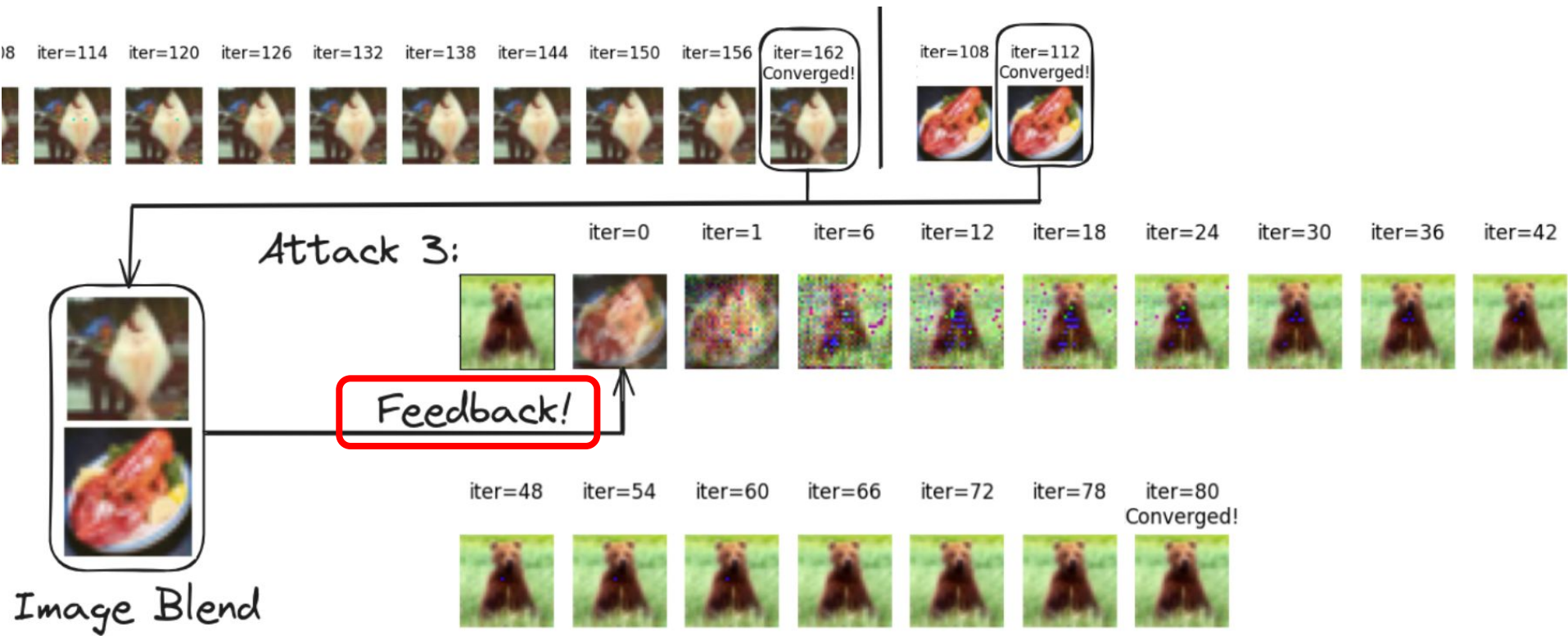
# Attack 2:

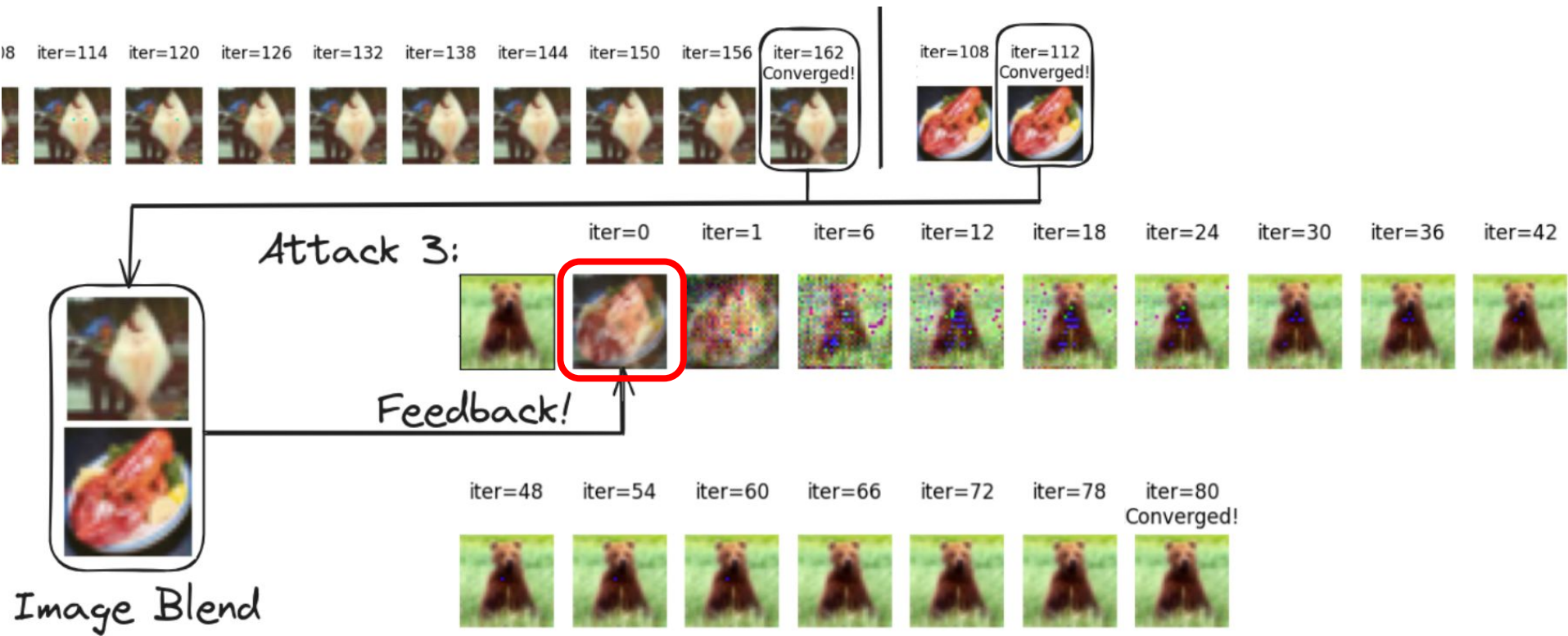


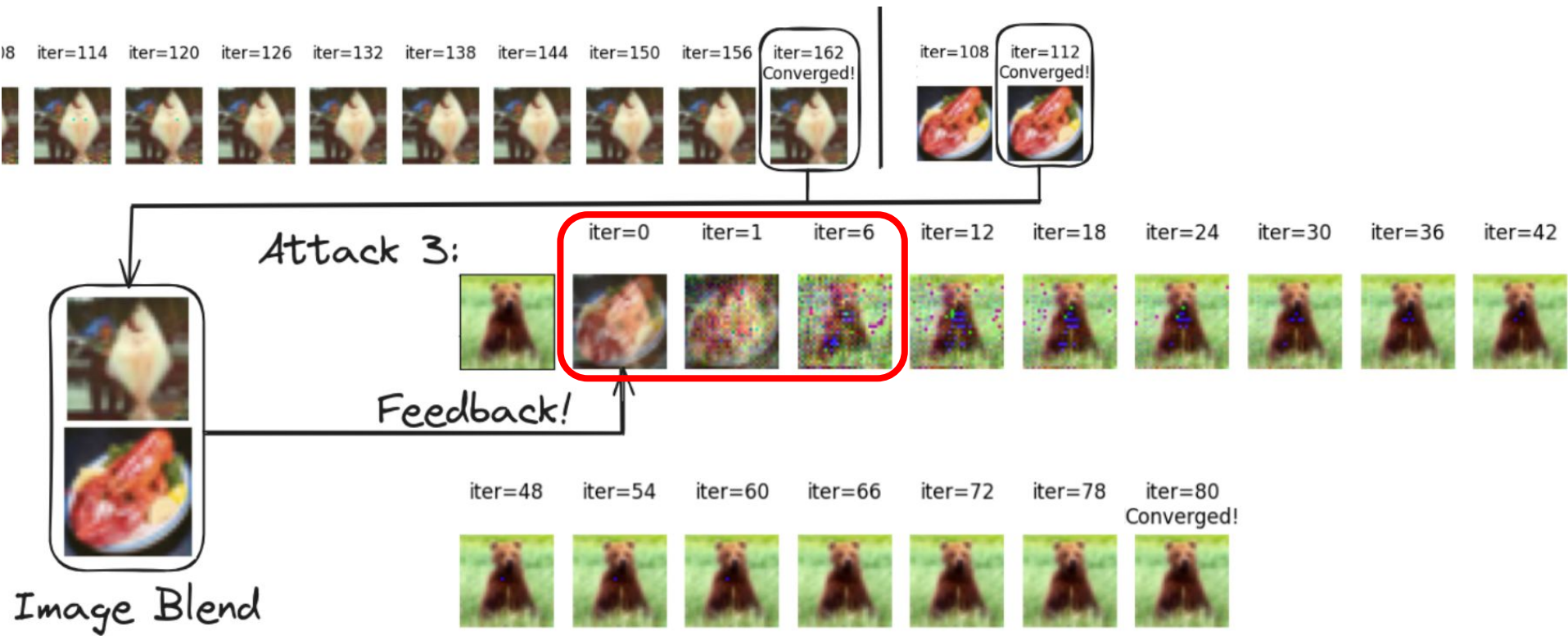


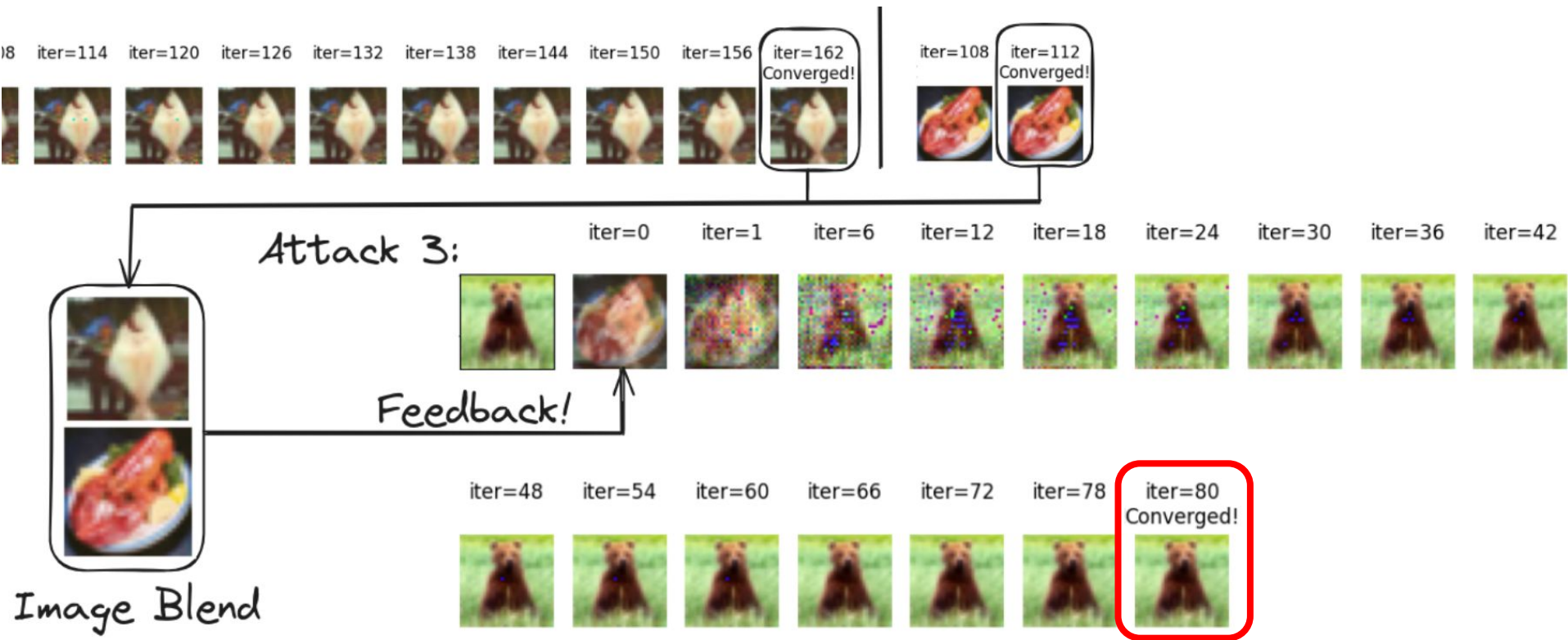


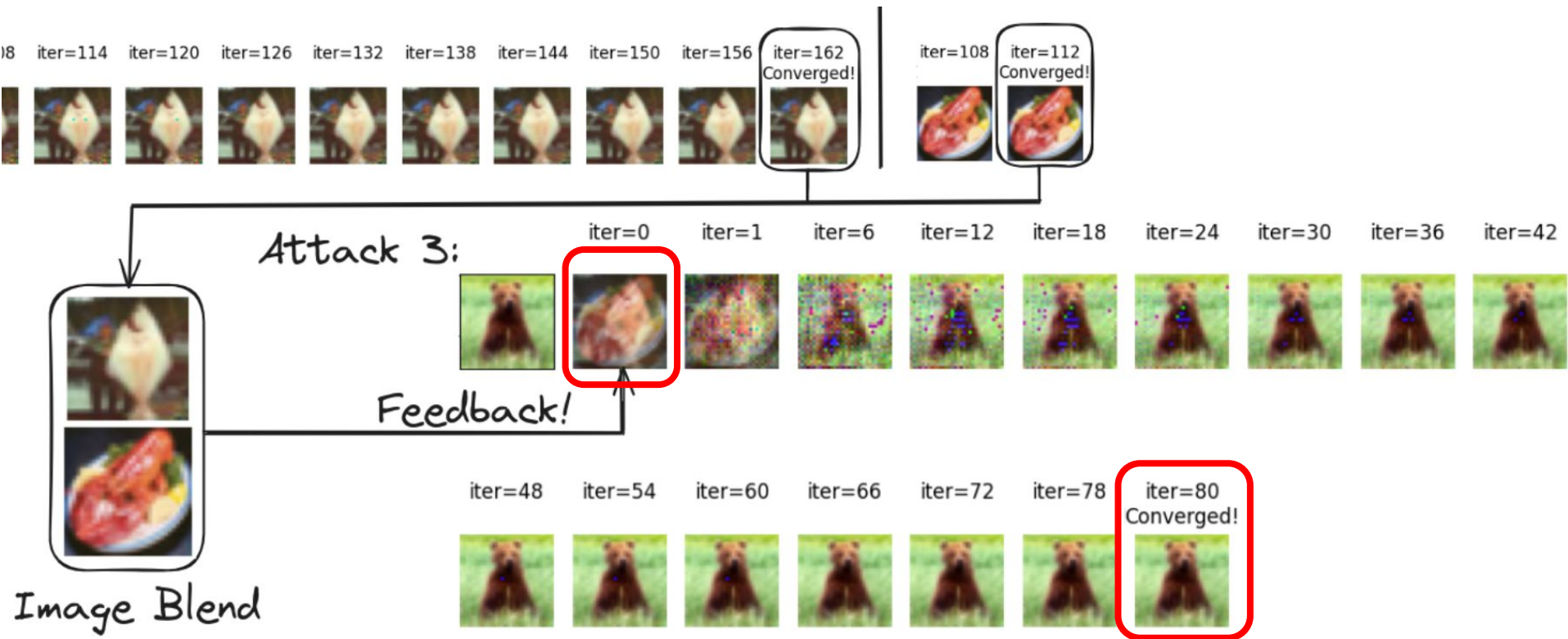


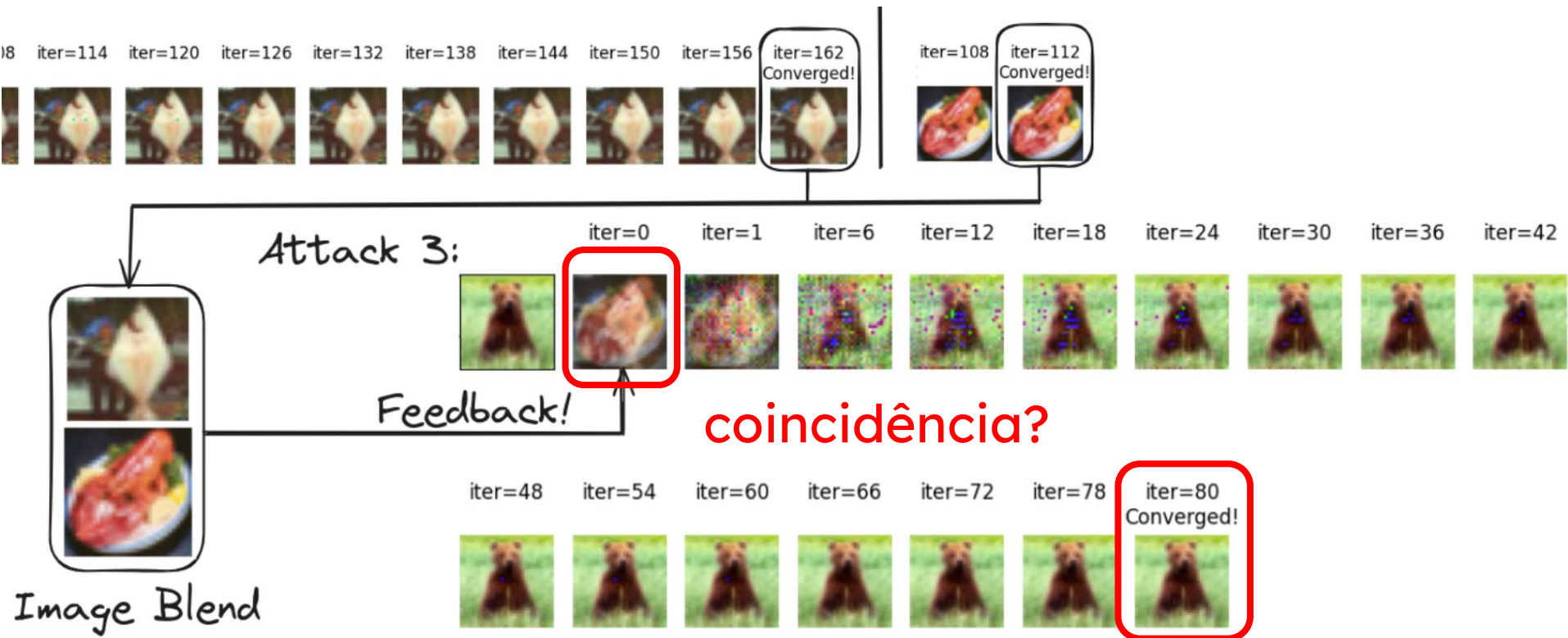






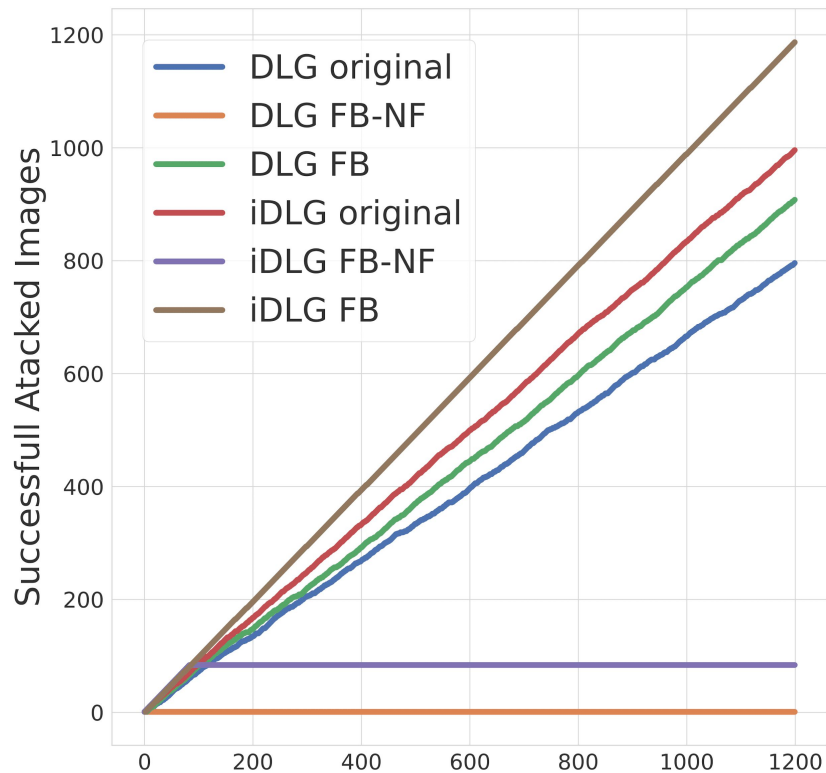




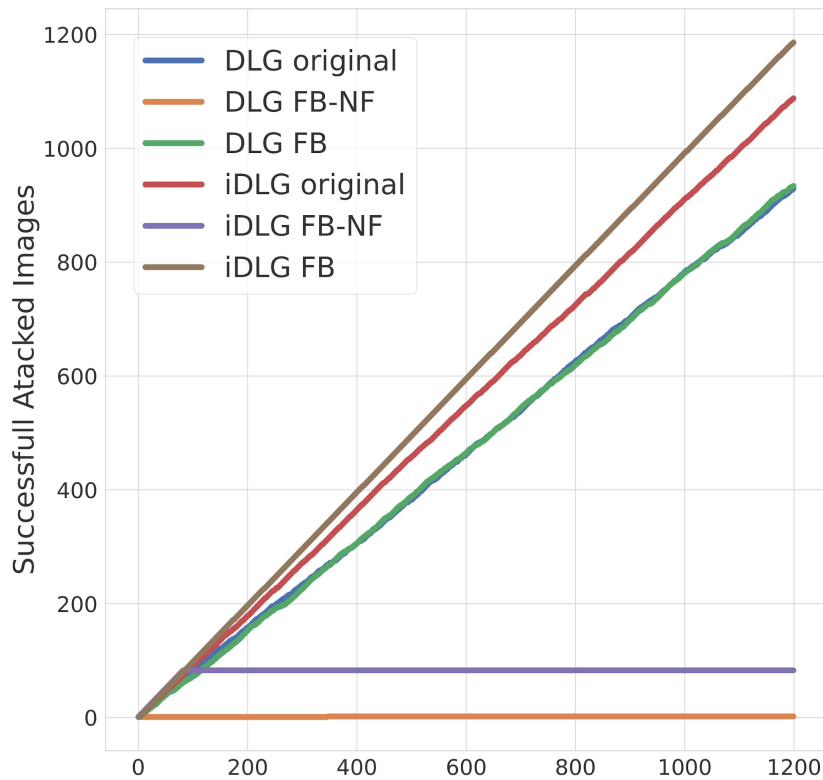




# Avaliação

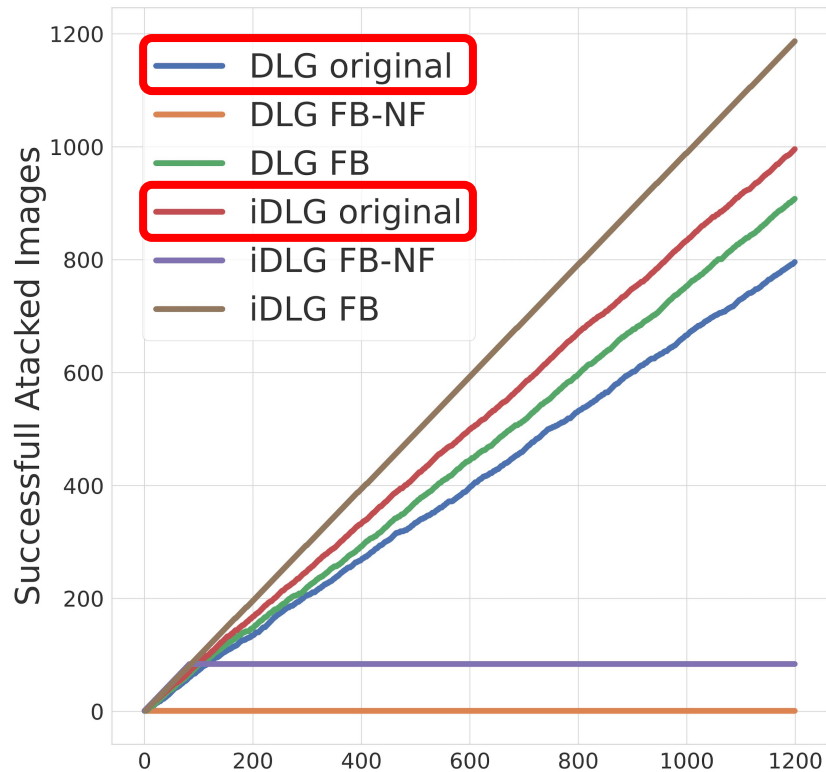


CIFAR100

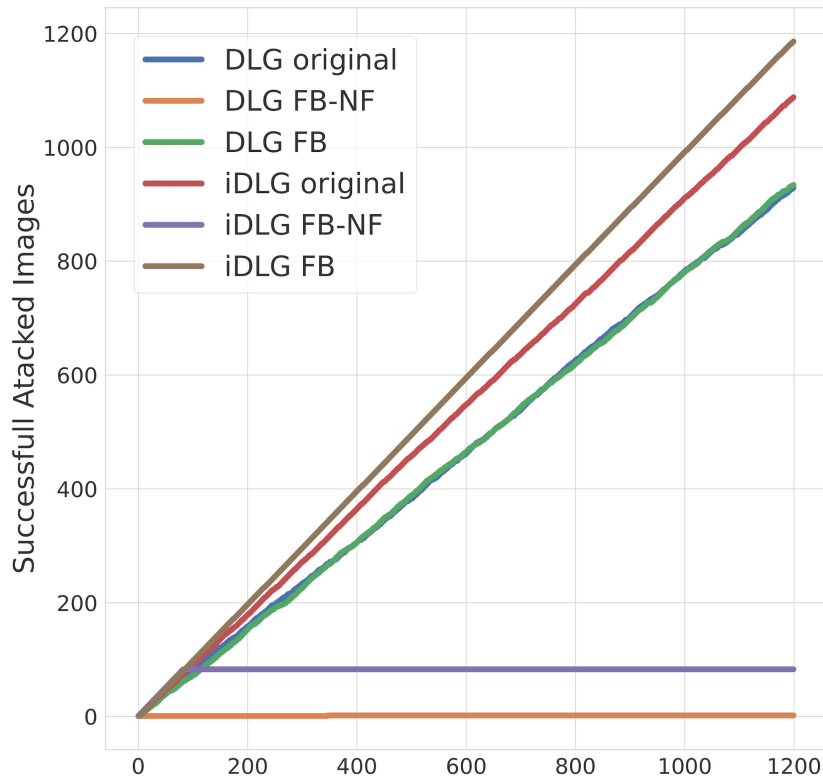


MNIST

# Avaliação

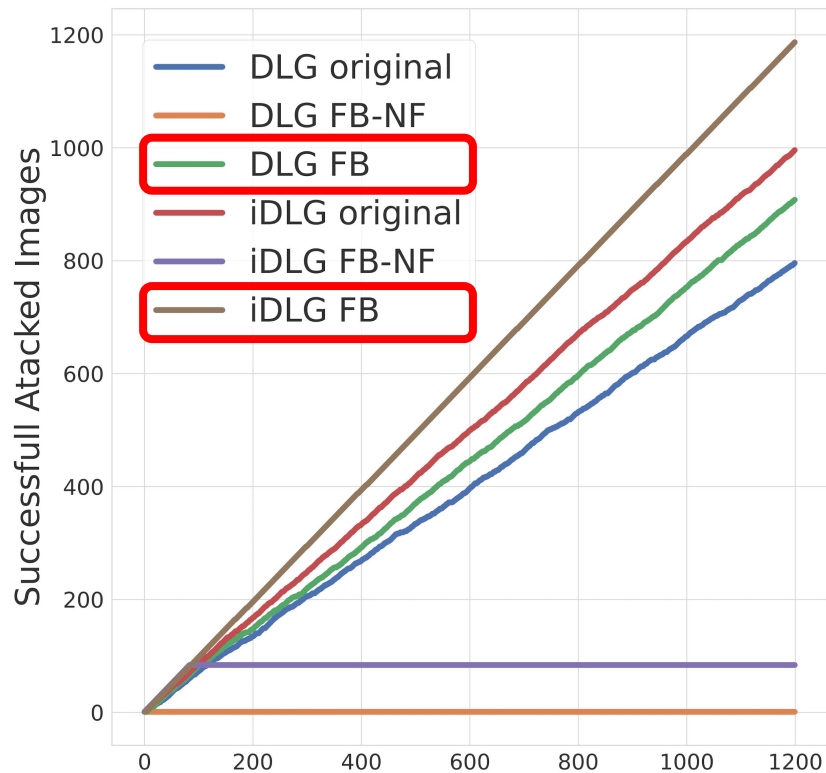


CIFAR100

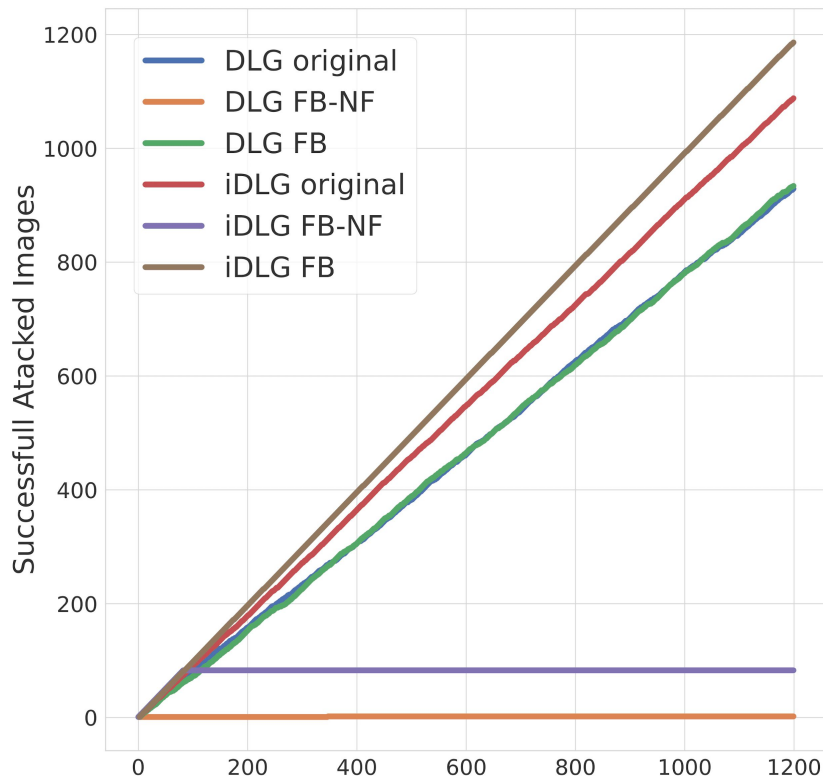


MNIST

# Avaliação

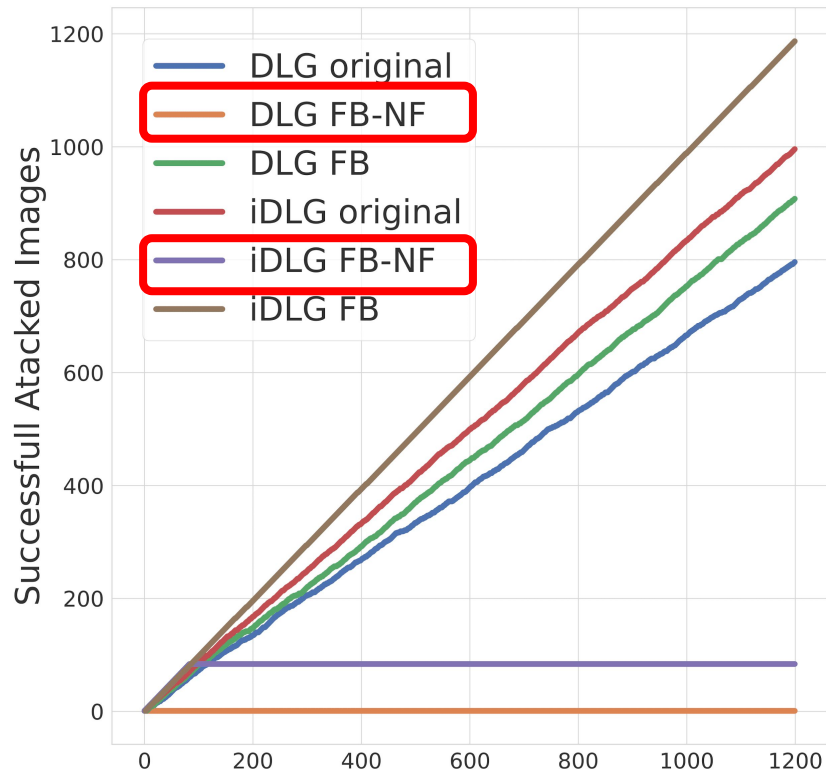


CIFAR100

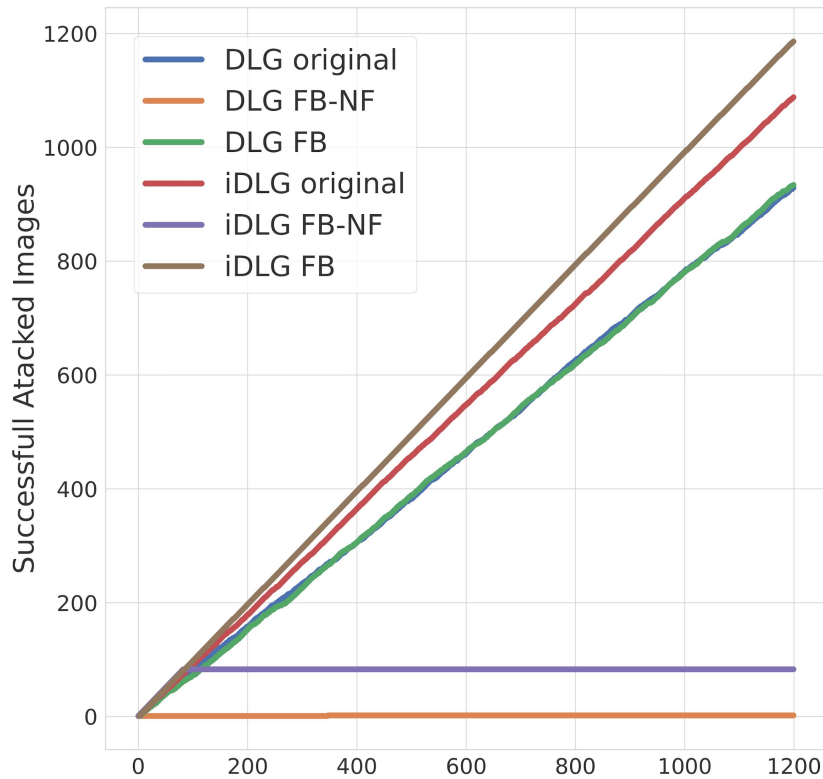


MNIST

# Avaliação

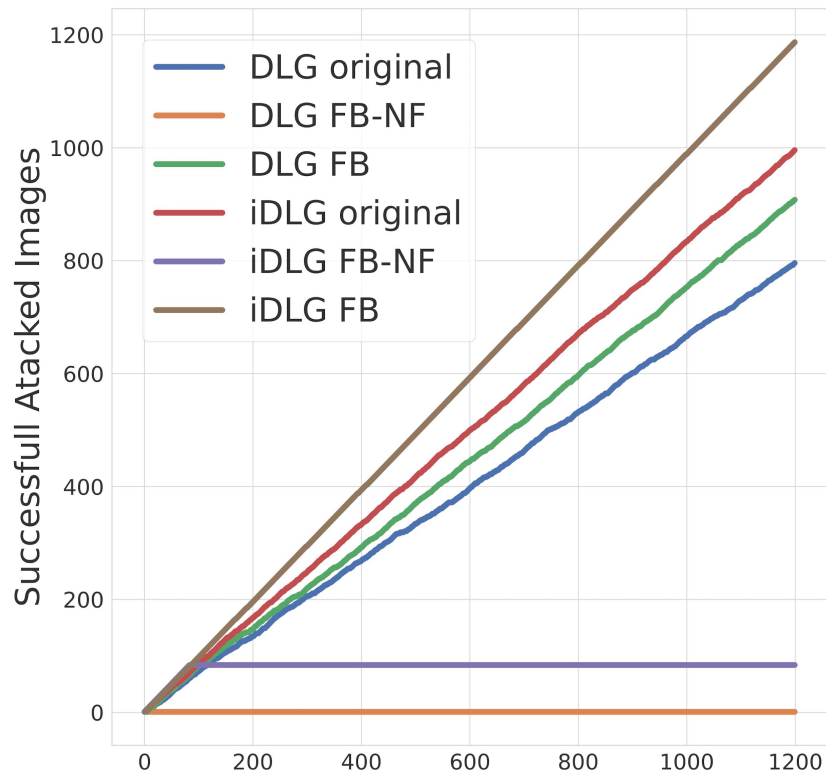


CIFAR100

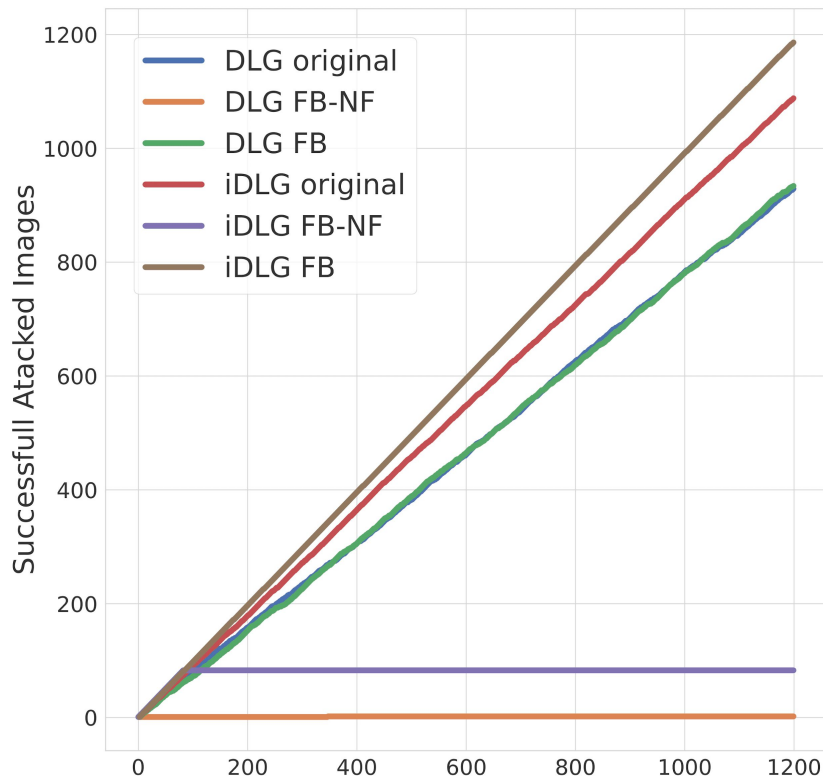


MNIST

# Avaliação

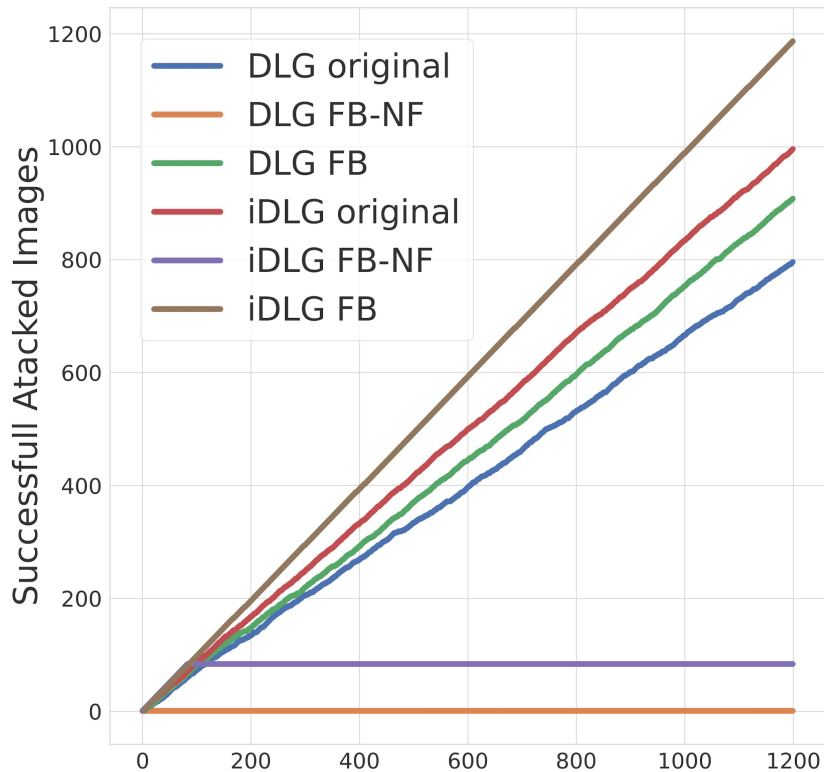


**CIFAR100**

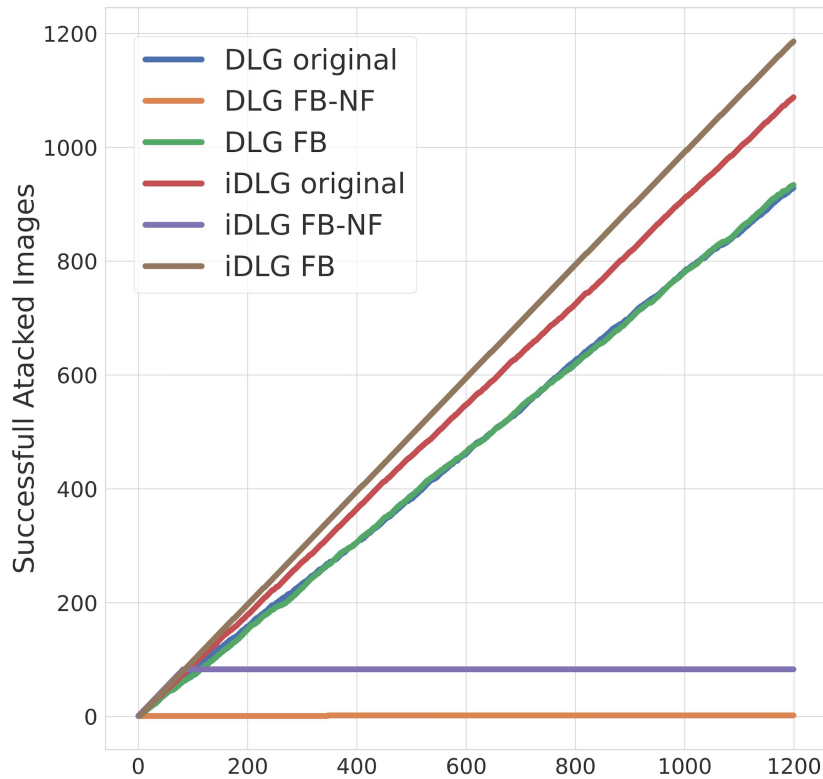


**MNIST**

# Avaliação

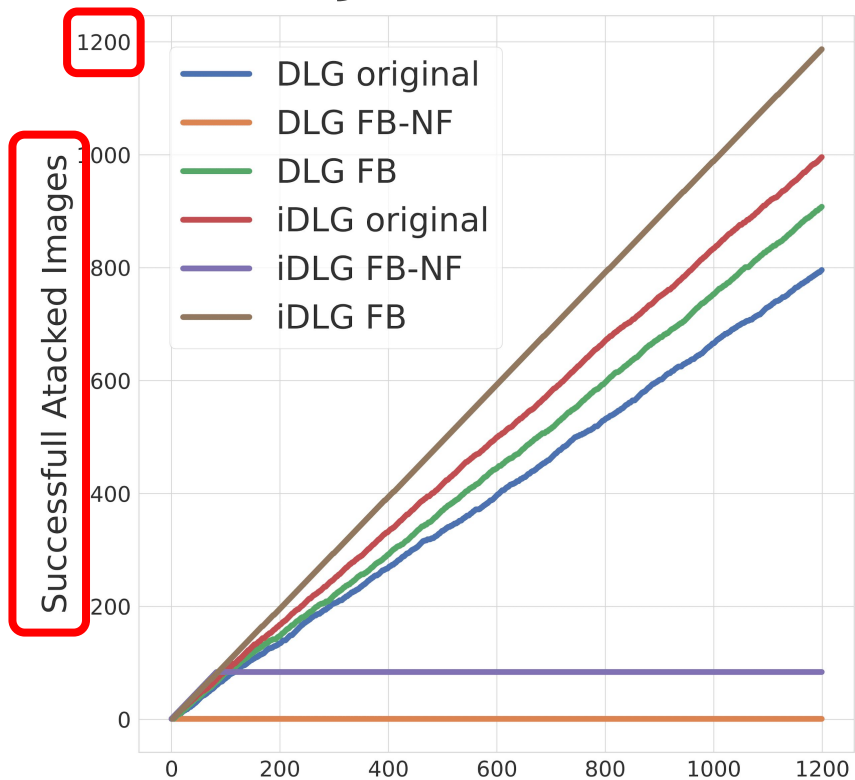


CIFAR100

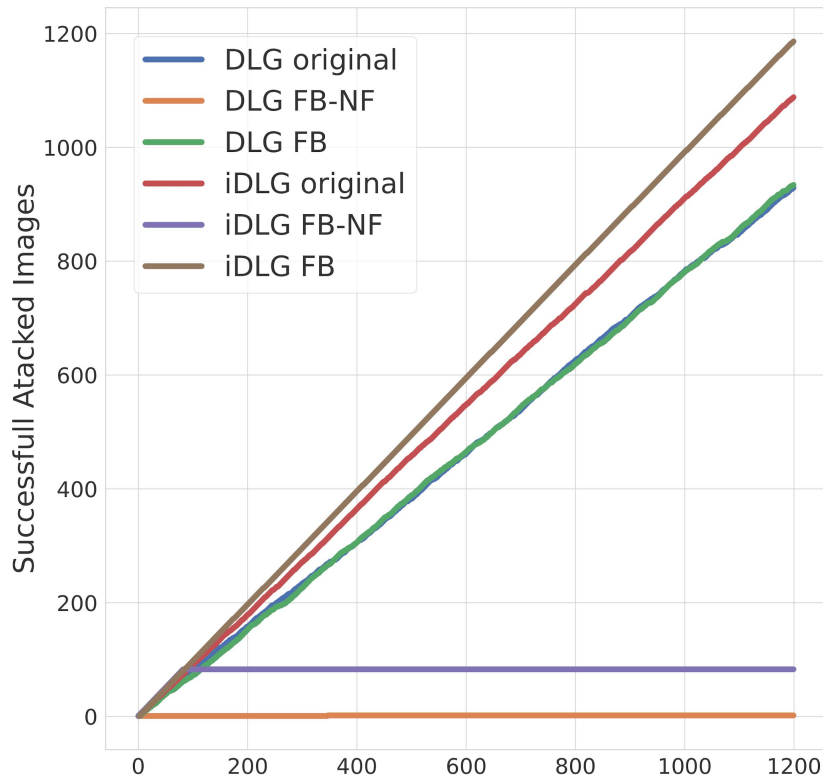


MNIST

# Avaliação

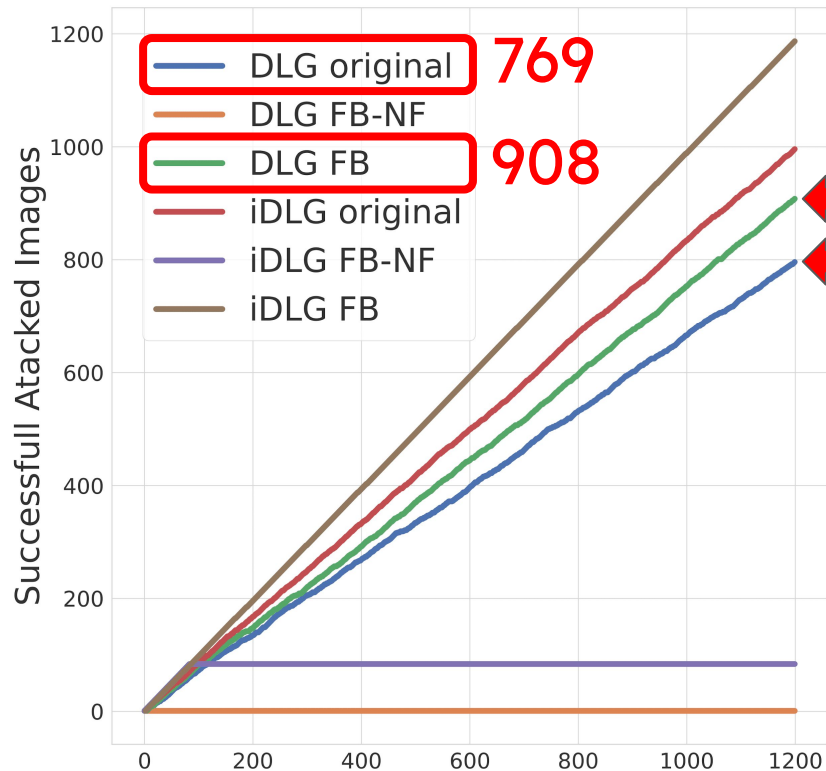


CIFAR100

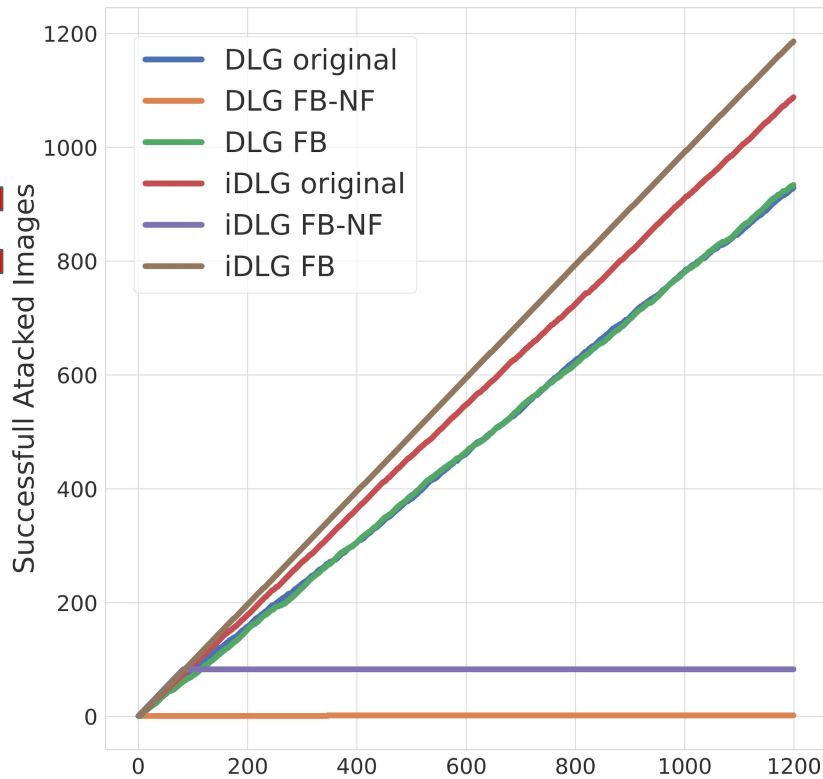


MNIST

# Avaliação +14.07%



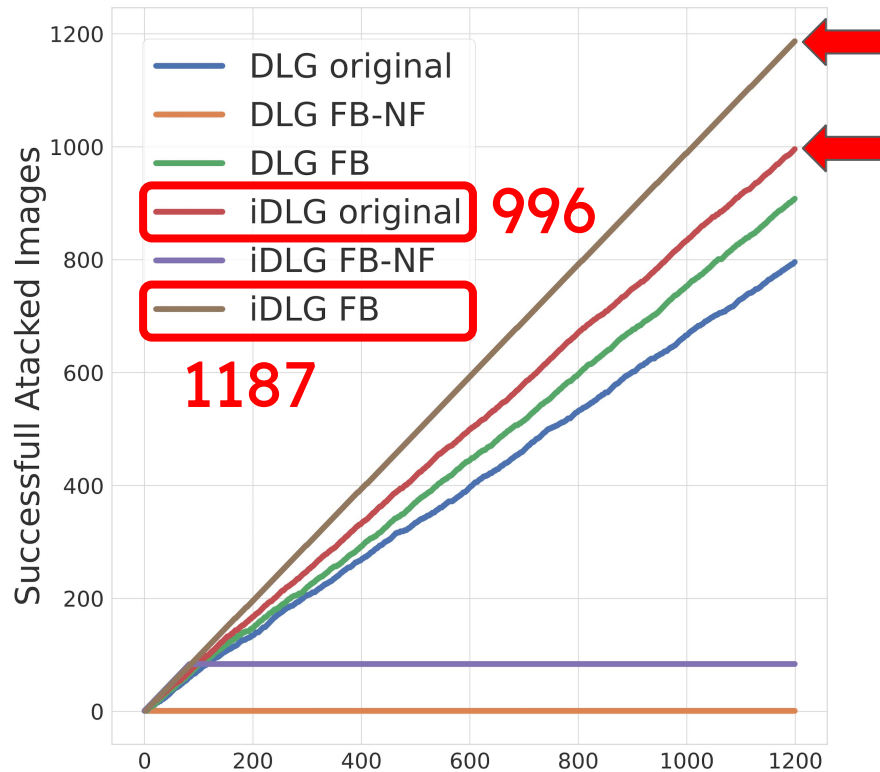
CIFAR100



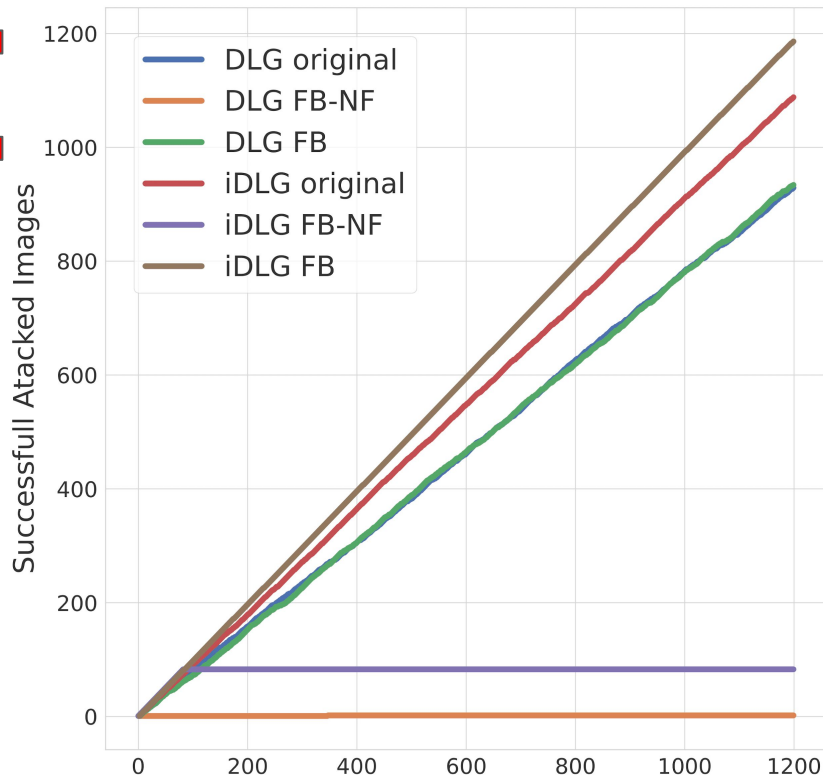
MNIST



# Avaliação +19.18%

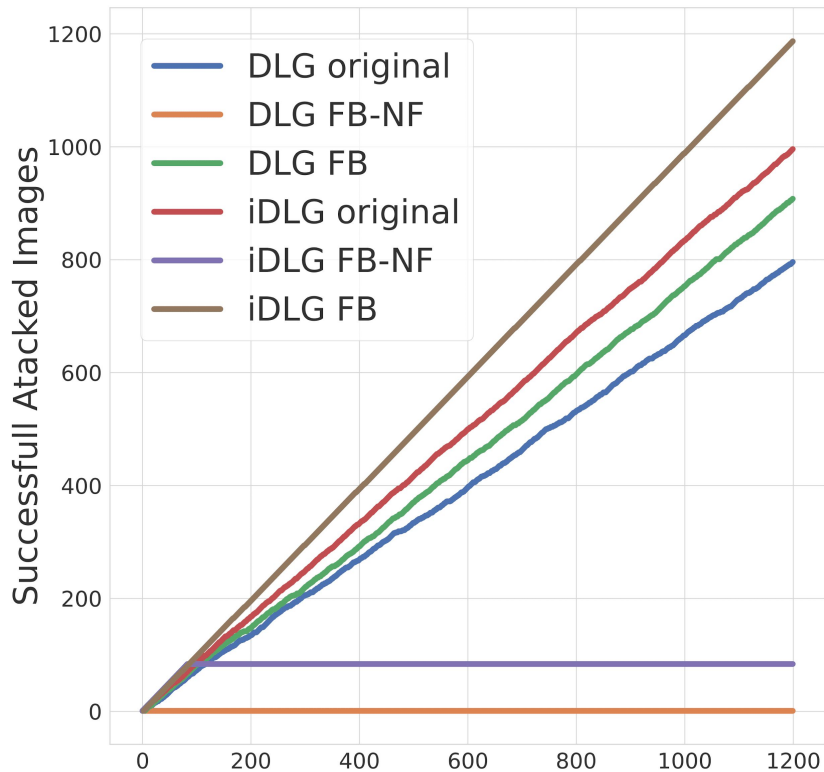


CIFAR100

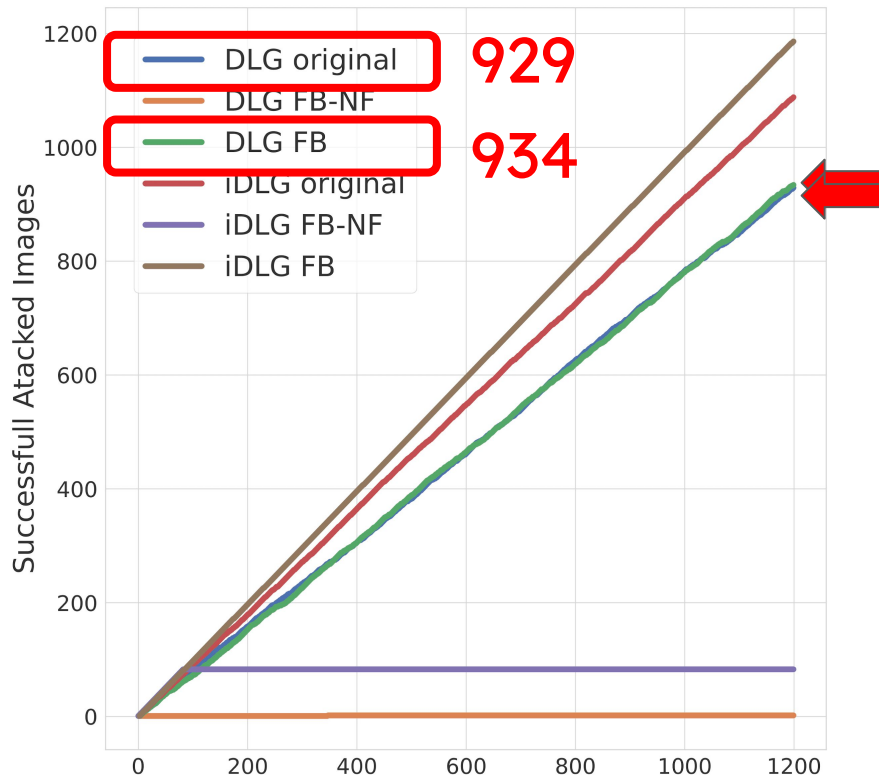


MNIST

# Avaliação +5 imagens

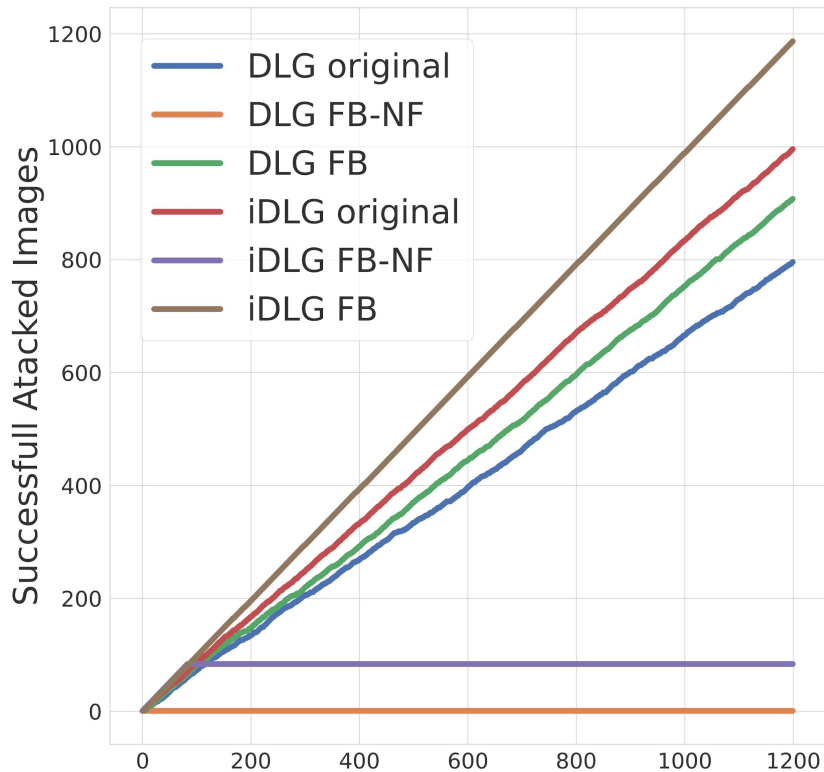


CIFAR100

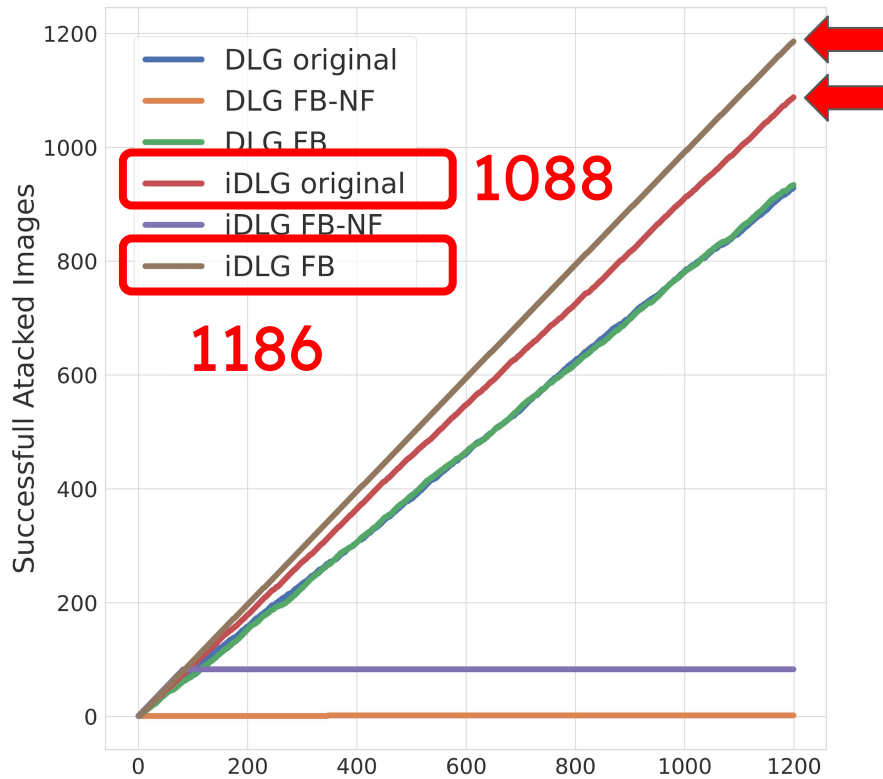


MNIST

# Avaliação +9.01%

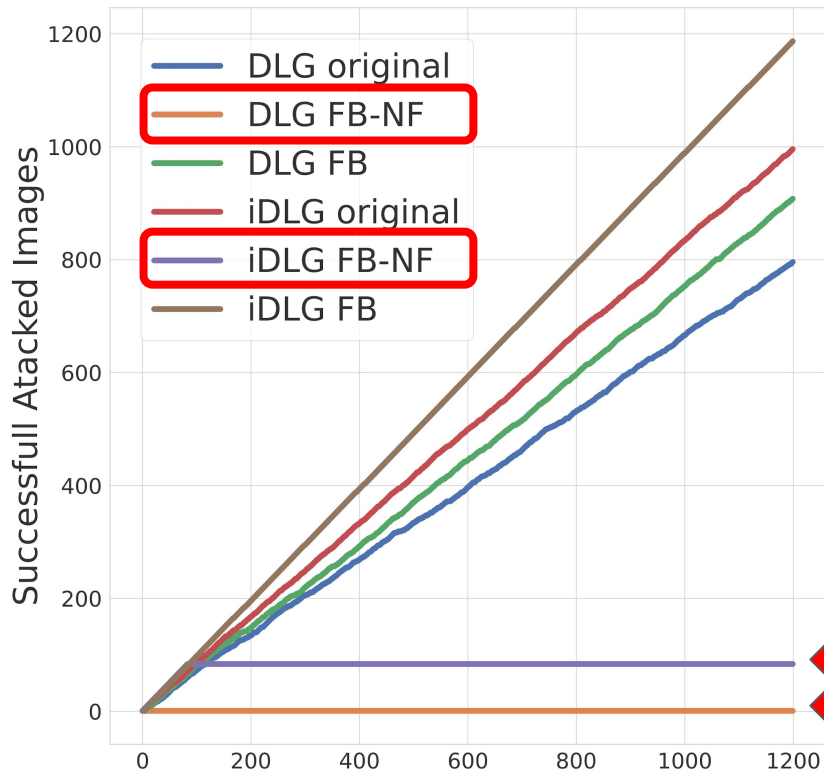


CIFAR100

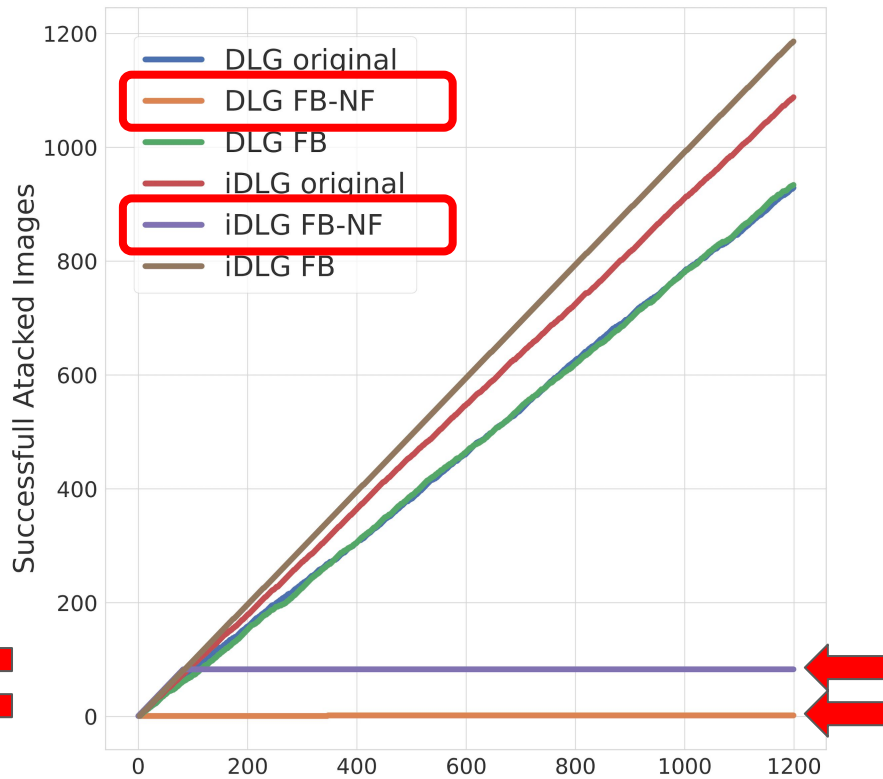


MNIST

# Avaliação

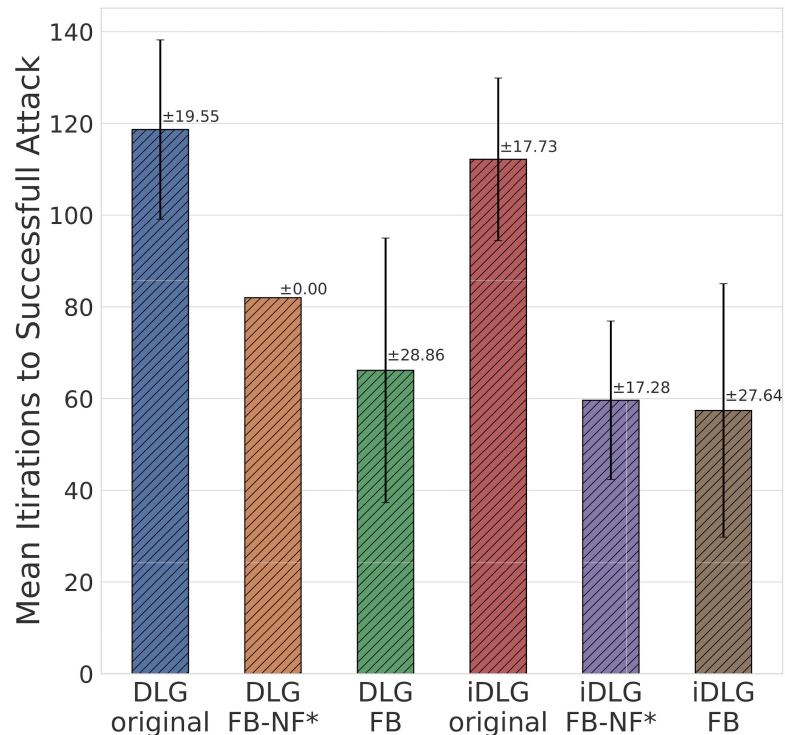


CIFAR100



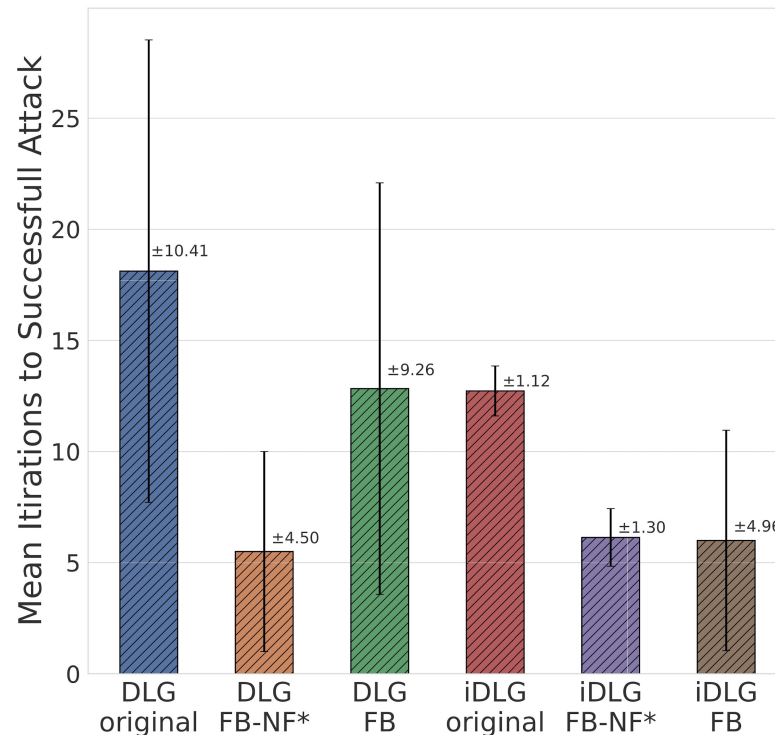
MNIST

# Avaliação



\*NF presents a reduced rate of attack success

## CIFAR100



\*NF presents a reduced rate of attack success

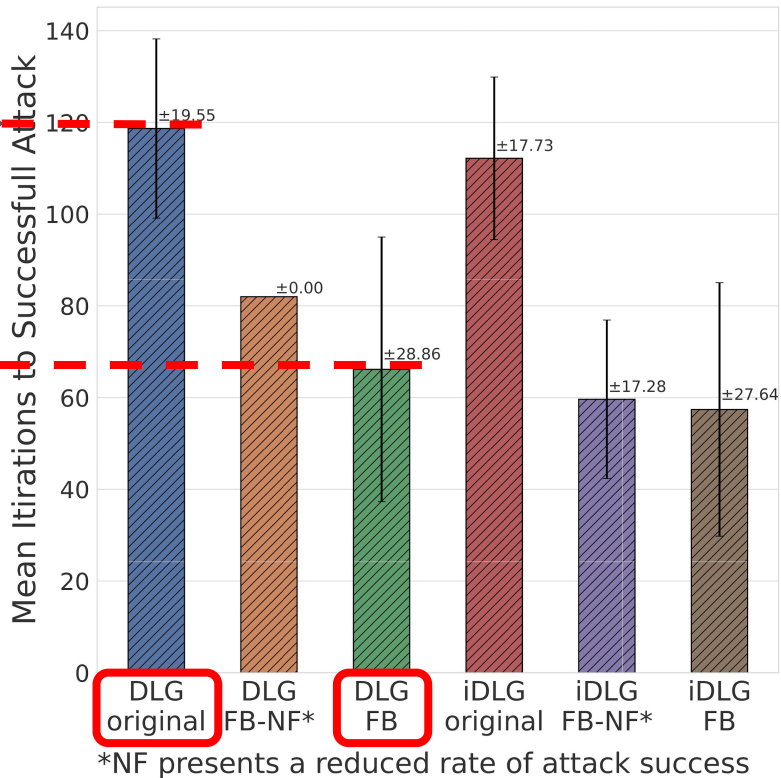
## MNIST

# Avaliação 44,26%

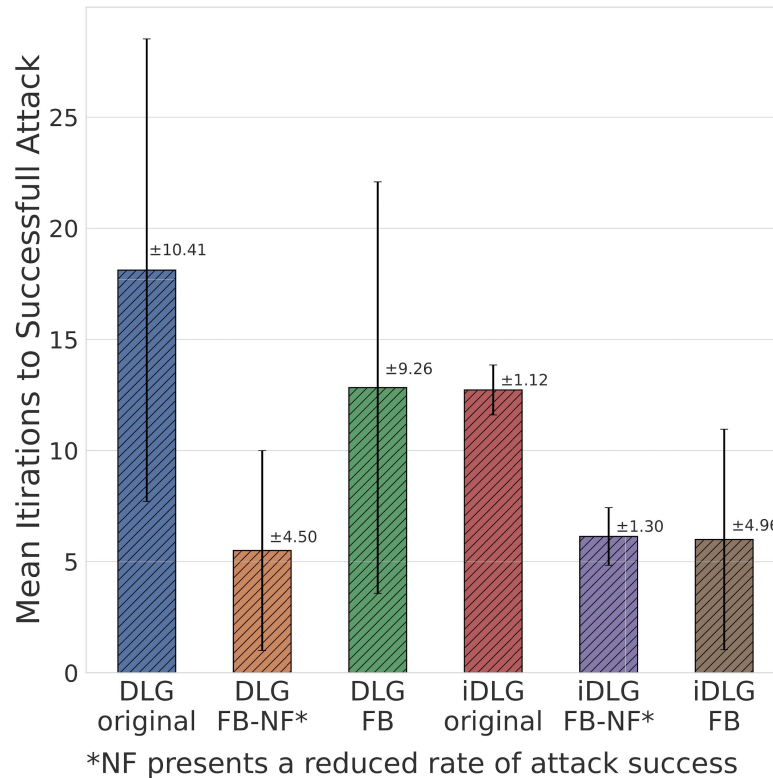
118



66

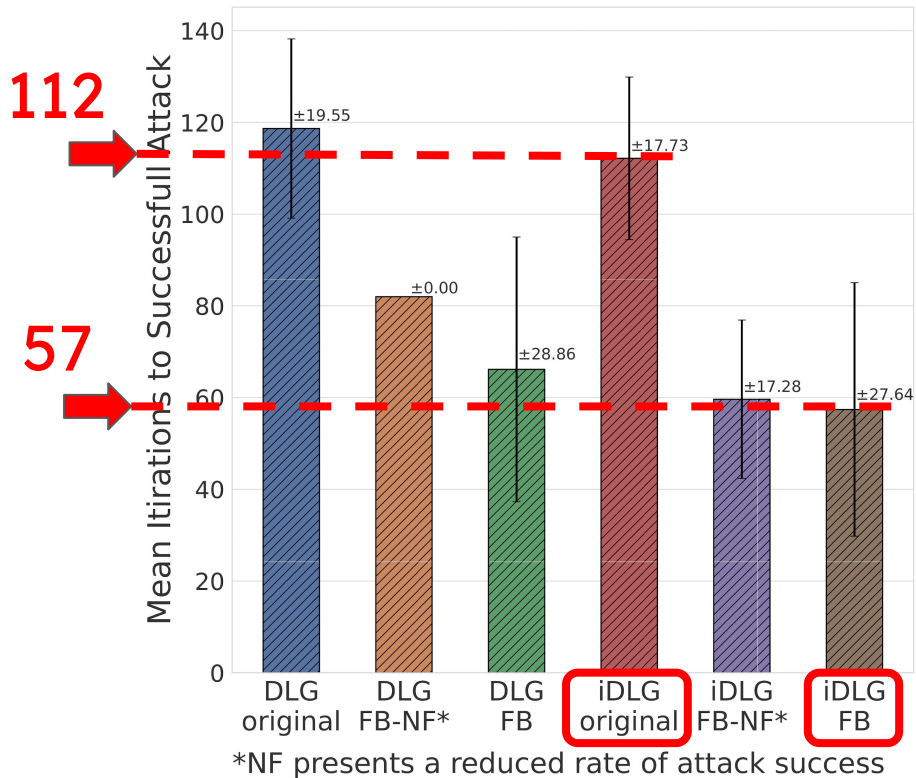


## CIFAR100

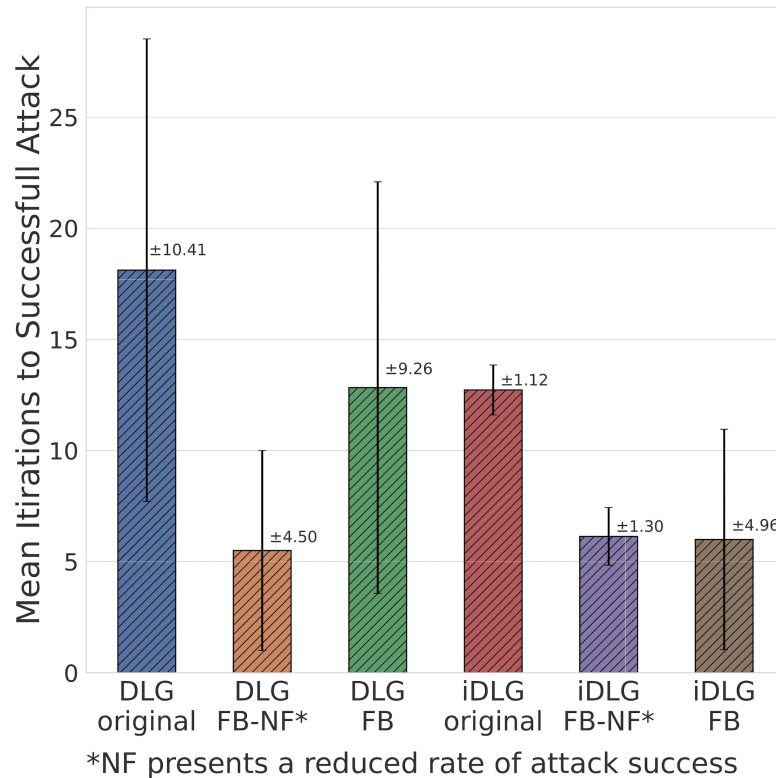


## MNIST

# Avaliação 48,82%

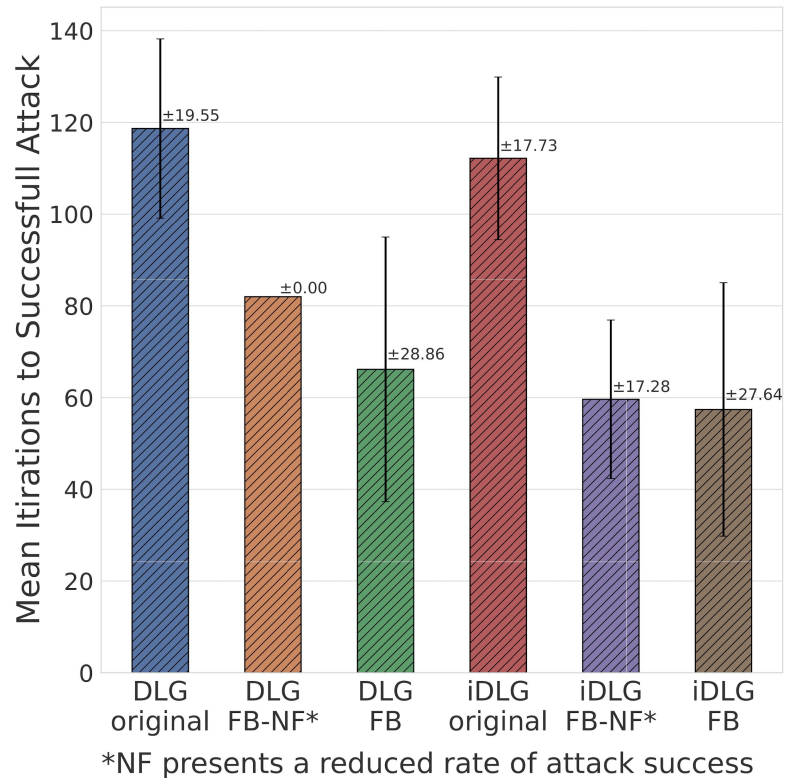


## CIFAR100

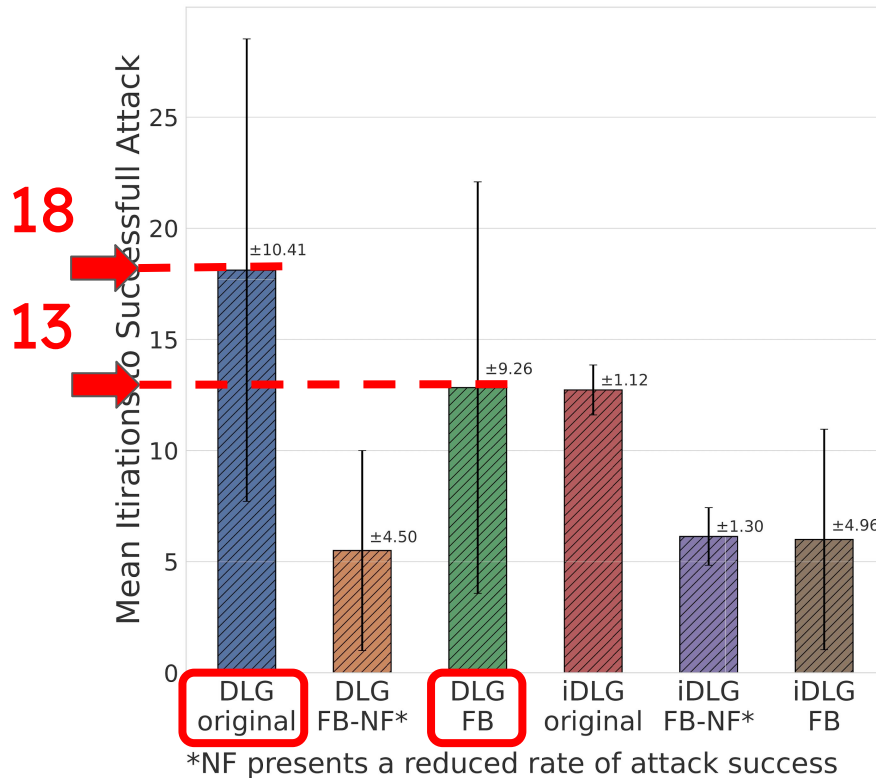


## MNIST

# Avaliação 29,19%



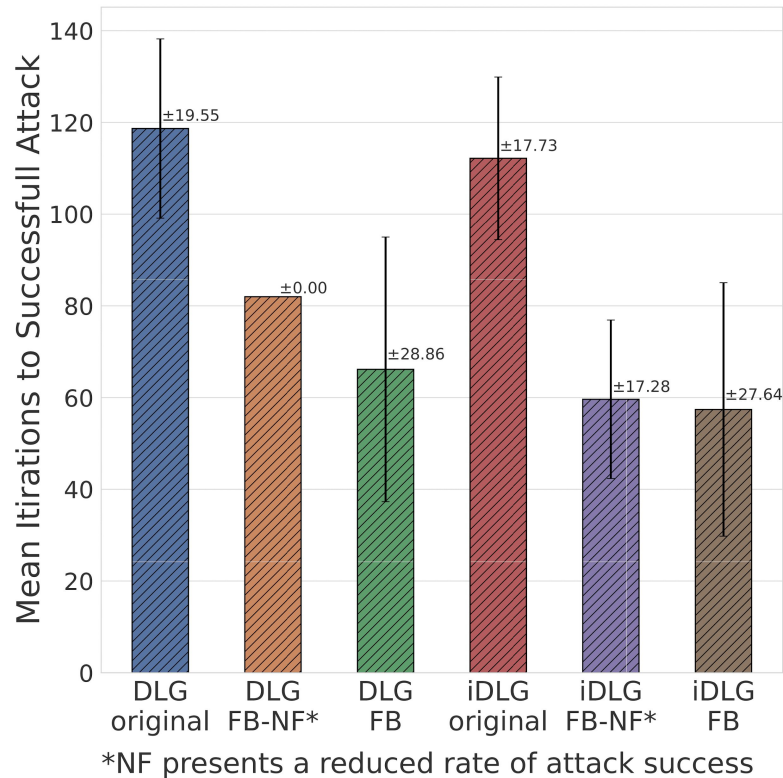
CIFAR100



MNIST



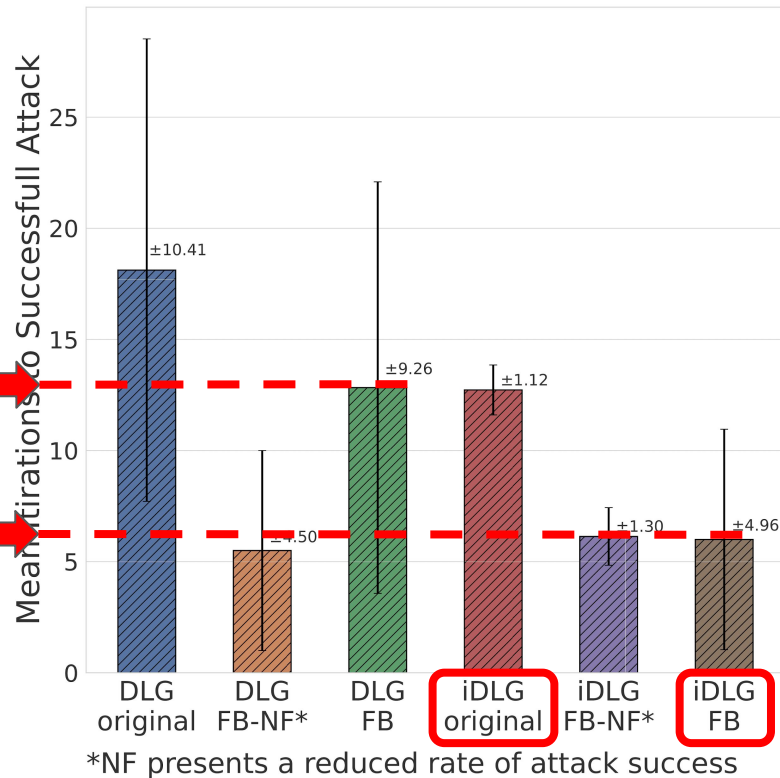
# Avaliação 52,59%



CIFAR100

13

6



MNIST

# Considerações finais

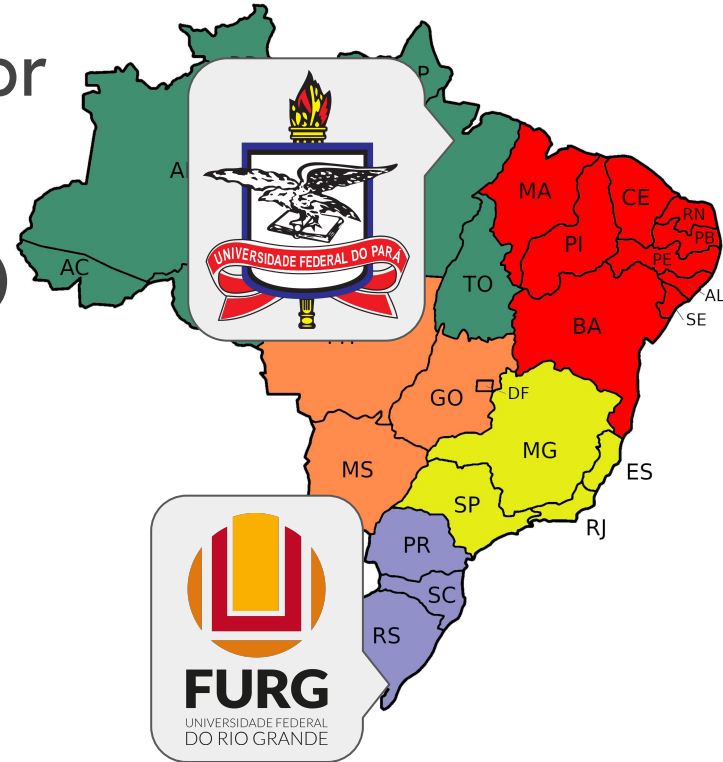
- **Sofisticação do ataque**
  - +taxa de sucesso
  - -custo computacional
- **Avanço contra 3 canais de cor**

# Trabalhos futuros

- Produzir mais métricas
- Analisar outros modelos
  - som, texto, vídeo
- Explorar defesas

# Obrigado!

- Luiz Antônio Leite
  - [luiz.freitas.leite@icen.ufpa.br](mailto:luiz.freitas.leite@icen.ufpa.br)
- Yuri Santo
- Prof. Bruno Dalmazo (FURG)
- Orientador:  
Prof. André Riker (UFPA)





# Patrocinadores do SBSeg 2024!

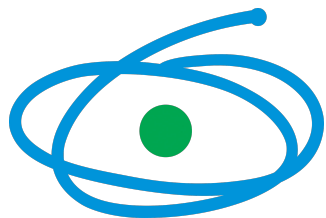
nie.br

egi.br

Google



Tempest



CAPES



SiDi



FAPESP



CNPq



C.E.S.A.R



zscaler™



BugHunt



FACULDADE  
IBPTech