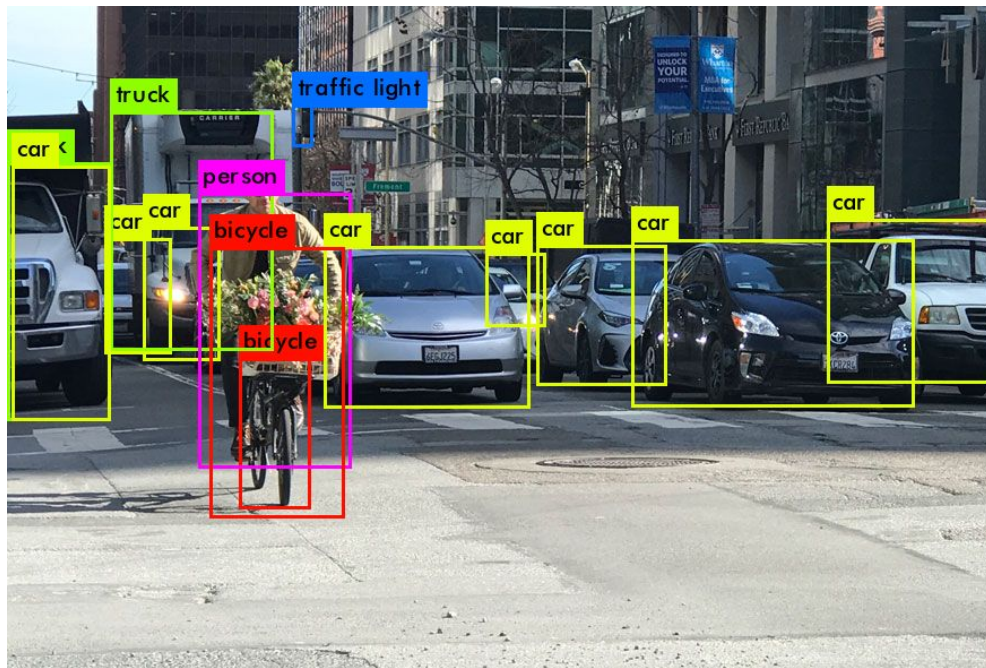# About me

➔ Final year undergraduate
➔ Deep Learning Researcher at FOR.ai
➔ Working on Adversarial Training and Robustness

FOR.ai

**Machine learning reached
state of the art on various tasks**

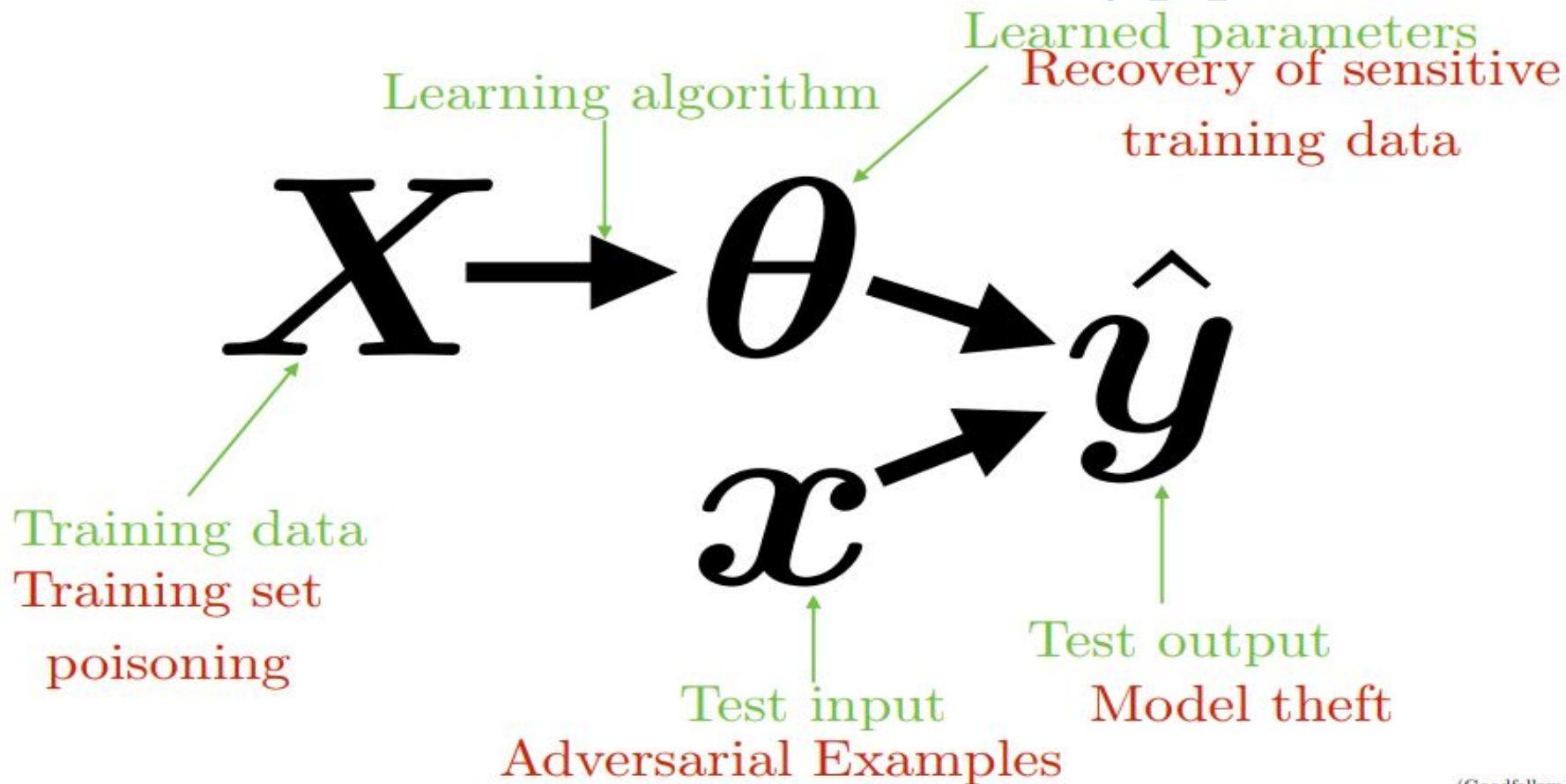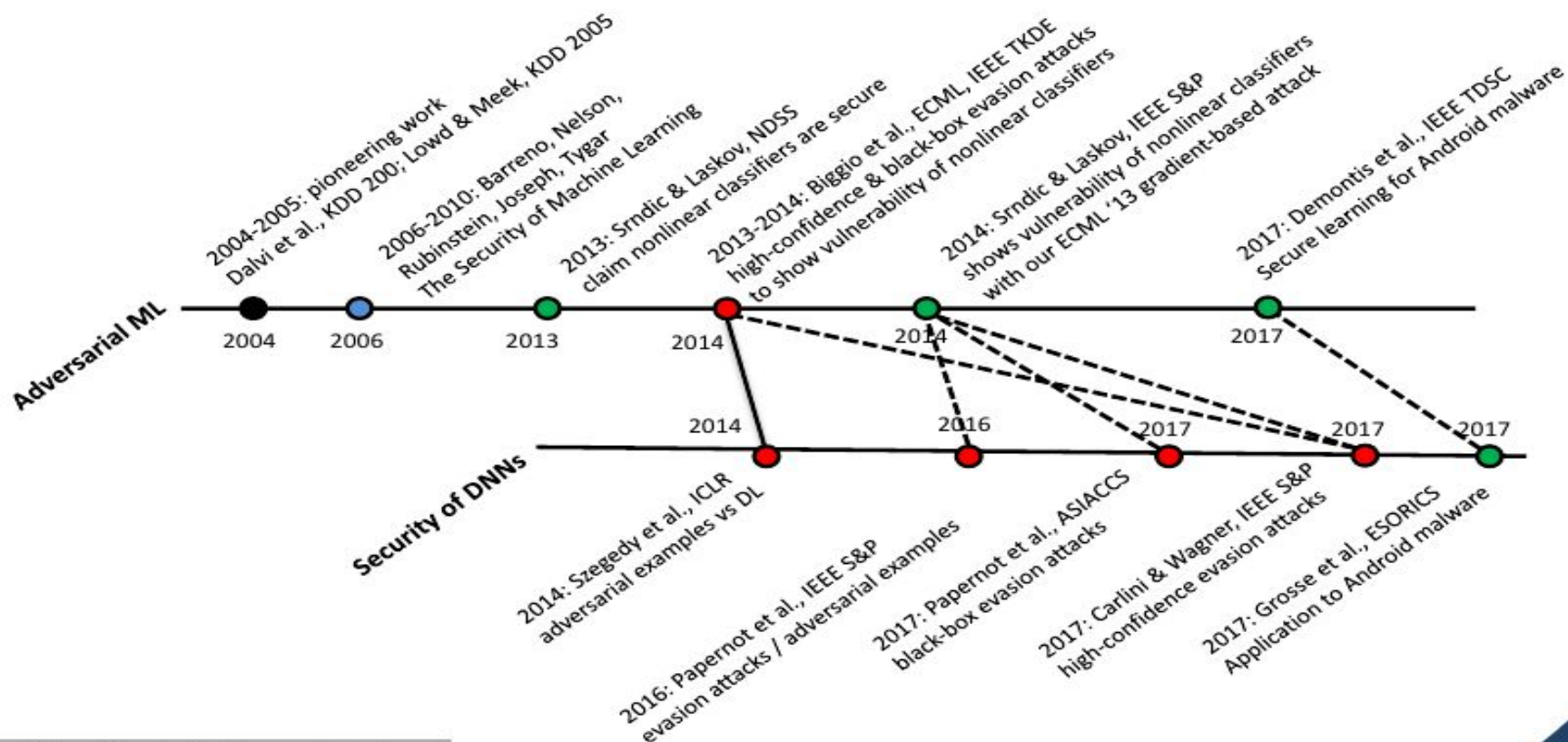# Good models make mistakes



"panda"
57.7% confidence

$+ \epsilon$

$=$

"gibbon"
99.3% confidence

# Attack on the machine learning pipeline

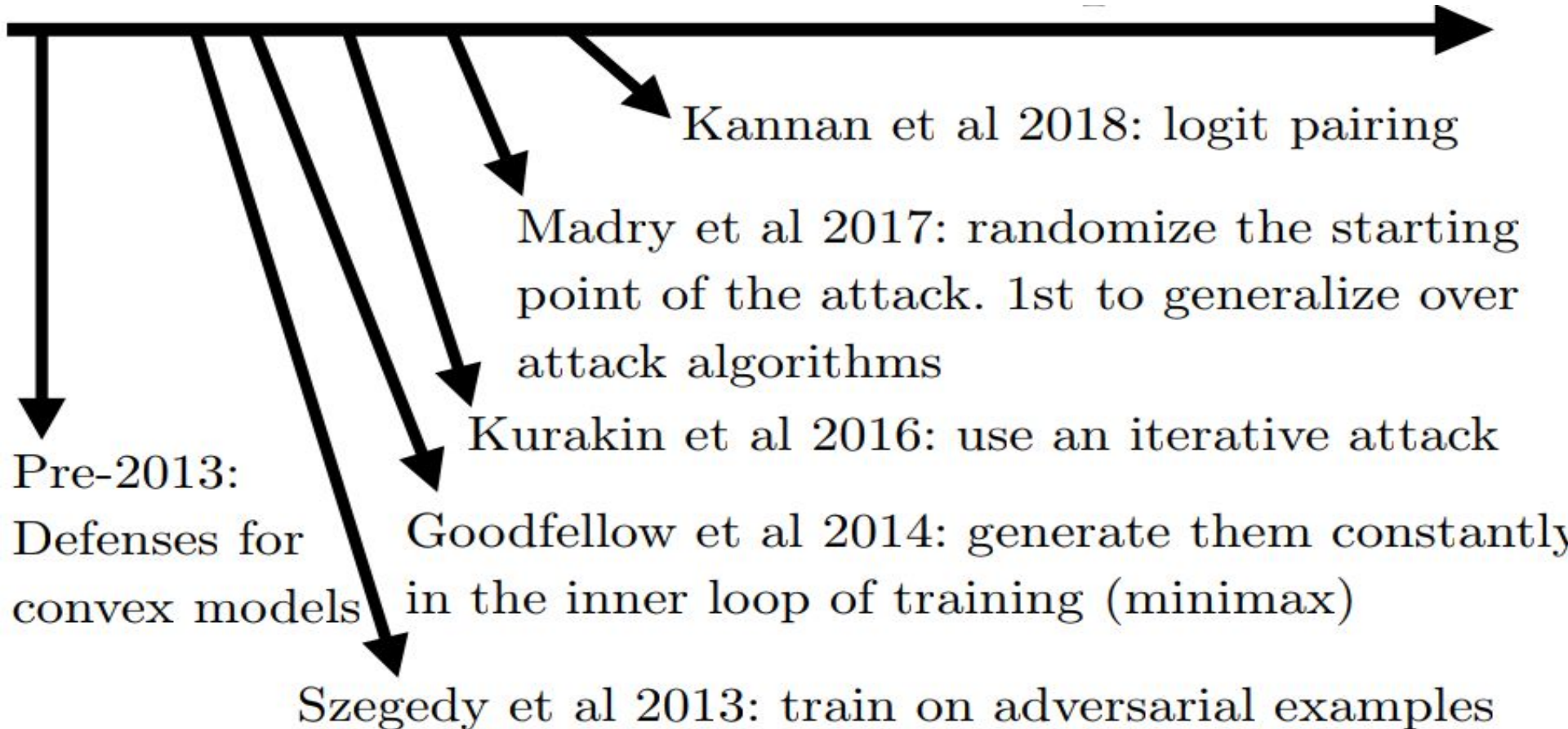

(Goodfellow 2018)

# Timeline of Learning Security

# What this means for us?

- Deep learning algorithms (Machine learning in general) are susceptible to attacks
- Use with caution in critical deployments
- Evaluate a model's adversarial resilience - not just accuracy/precision/recall
- Spend effort to make model robust to tempering

# Defending the machines

- Distillation (Train model 2x, feed first DNN output logits into second DNN input layer)
- Train models with adversarial samples i.e ironing out imperfect knowledge learnt in the model)
- Special regularization methods/loss functions (simulating adversarial content during training)

# Timeline of Defences



Kannan et al 2018: logit pairing

Madry et al 2017: randomize the starting point of the attack. 1st to generalize over attack algorithms

Kurakin et al 2016: use an iterative attack

Pre-2013:
Defenses for convex models

Goodfellow et al 2014: generate them constantly in the inner loop of training (minimax)

Szegedy et al 2013: train on adversarial examples

(Goodfellow 2018)

# Thank you
# Questions?

Twitter: @divyam3897
Linkedin: https://www.linkedin.com/in/dmadaan3897
Github: @divyam3897