

Machine Learning

in the presence of adversaries

Andrew Paverd and Mika Juuti
(joint work with N. Asokan, Jian Liu and Samuel Marchal)

About Us

Aalto Department of Computer Science

<http://cs.aalto.fi/>

- Top-100 in CS world rankings
- AI, algorithms, security & privacy

Secure Systems Group (Prof. Asokan)

http://cs.aalto.fi/secure_systems

- 10 senior researchers, 5-10 MSc students

Here today:

- Andrew Paverd: PhD from Oxford 2016
- Mika Juuti: M.Sc. (Tech) from Aalto 2015

<https://ajpaverd.org/>

<https://research.aalto.fi/portal/mika.juuti.html>

Machine Learning is ubiquitous

The ML market size is expected to grow by **44% annually** over next five years
In 2016, companies invested up to **\$9 Billion** in AI-based startups



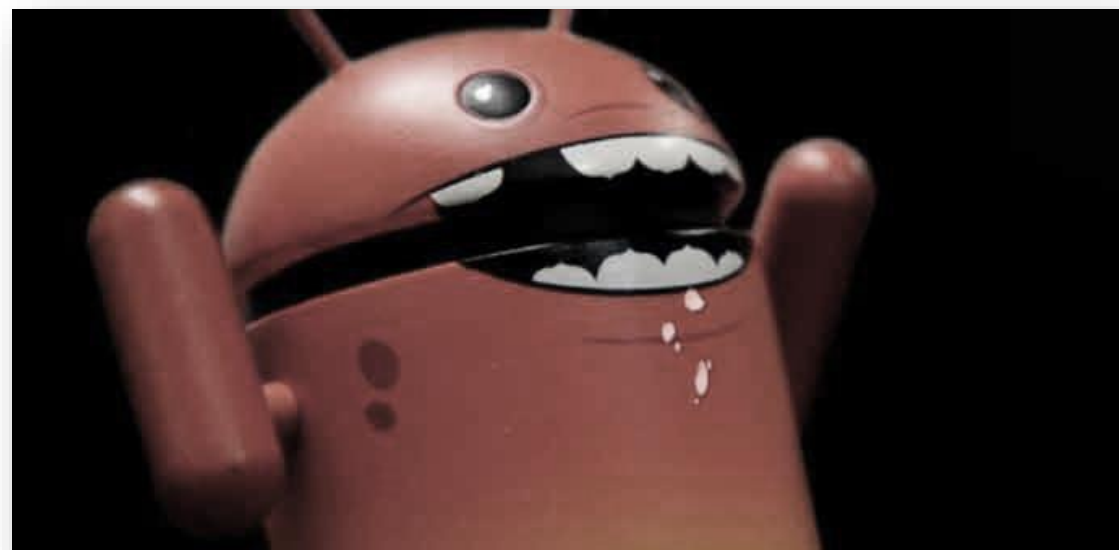
[1] <http://www.marketsandmarkets.com/PressReleases/machine-learning.asp>

[2] McKinsey Global Institute, "Artificial Intelligence: The Next Digital Frontier?"

Machine Learning *for* security/privacy



Access Control



Malware / Intrusion Detection

Security & privacy of machine learning



Which class is this?
School bus



Which class is this?
Ostrich



Which class is this?

Panda

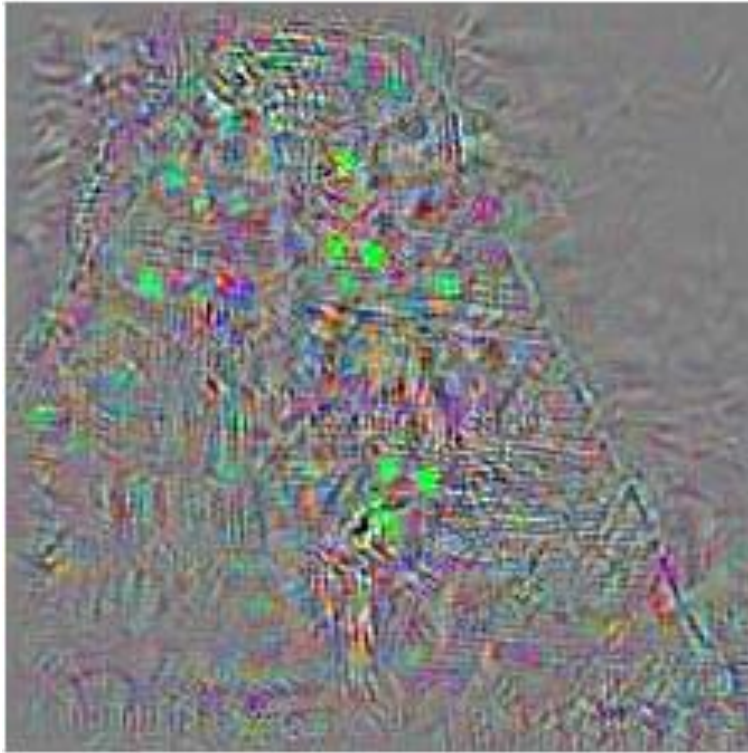


Which class is this?

Gibbon



Which class is this?
Building



Which class is this?
Ostrich



Which class is this?
Cat



Which class is this?
Desktop computer

Athalye et al. "Synthesizing Robust Adversarial Examples"

<https://blog.openai.com/robust-adversarial-inputs/>

DolphinAttack: Inaudible Voice command

Guoming Zhang Chen Yan Xiaoyu Ji

Tianchen Zhang Taimin Zhang Wenyan Xu

Zhejiang University

ACM CCS 2017



Target



Softmax

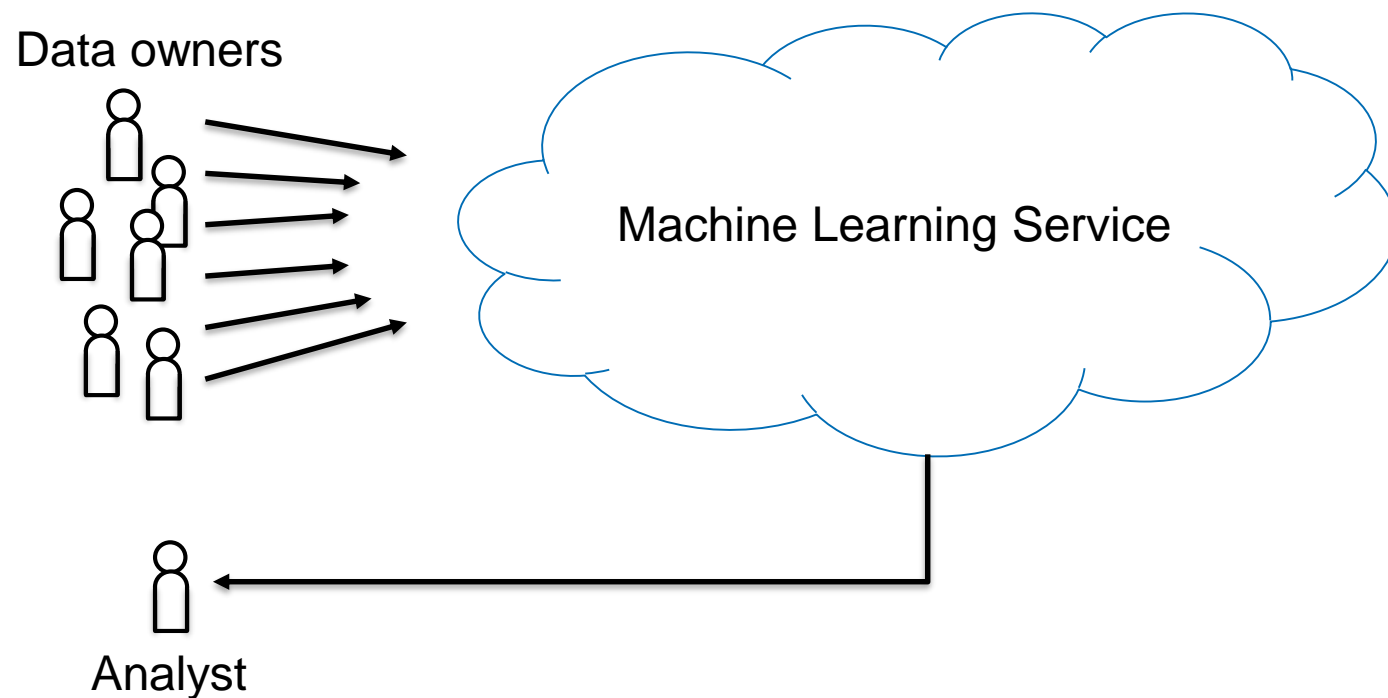


MLP

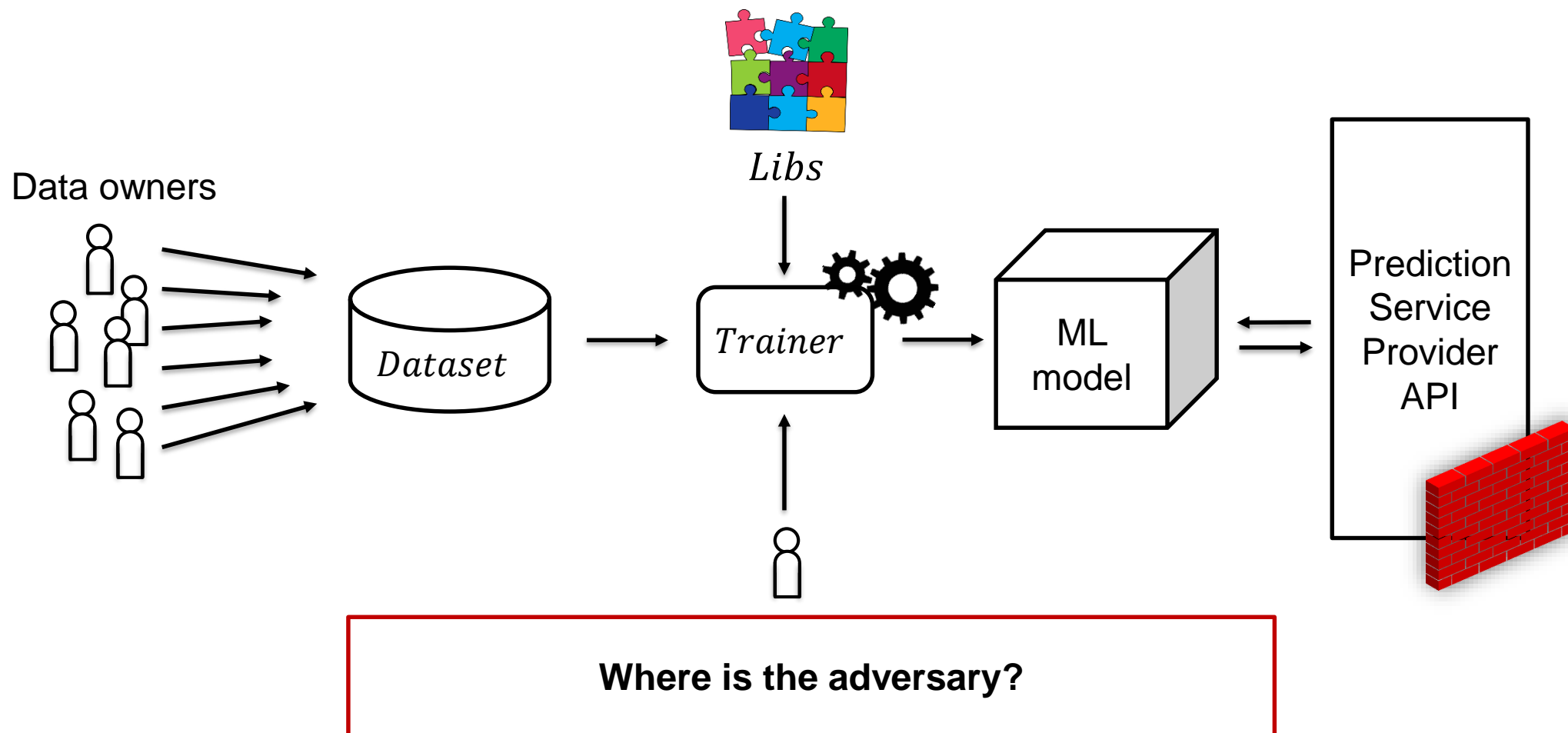


DAE

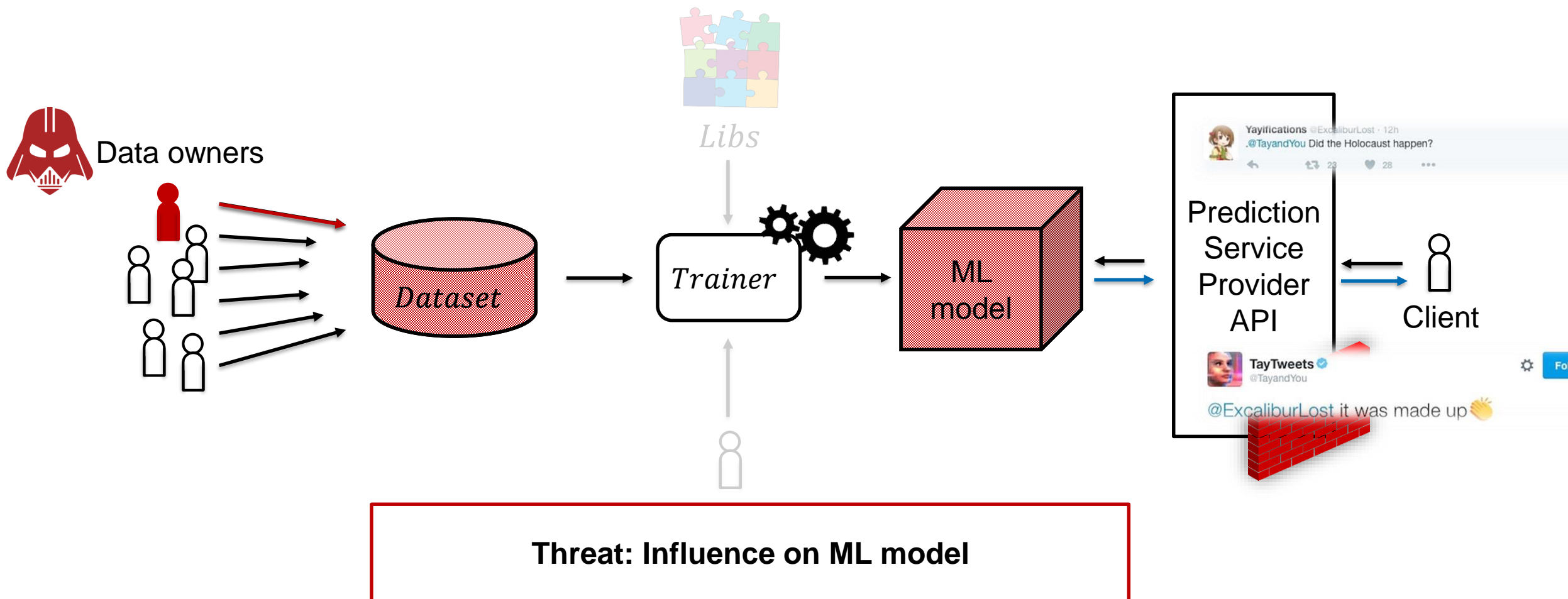
A Basic Machine Learning pipeline



A more realistic Machine Learning pipeline

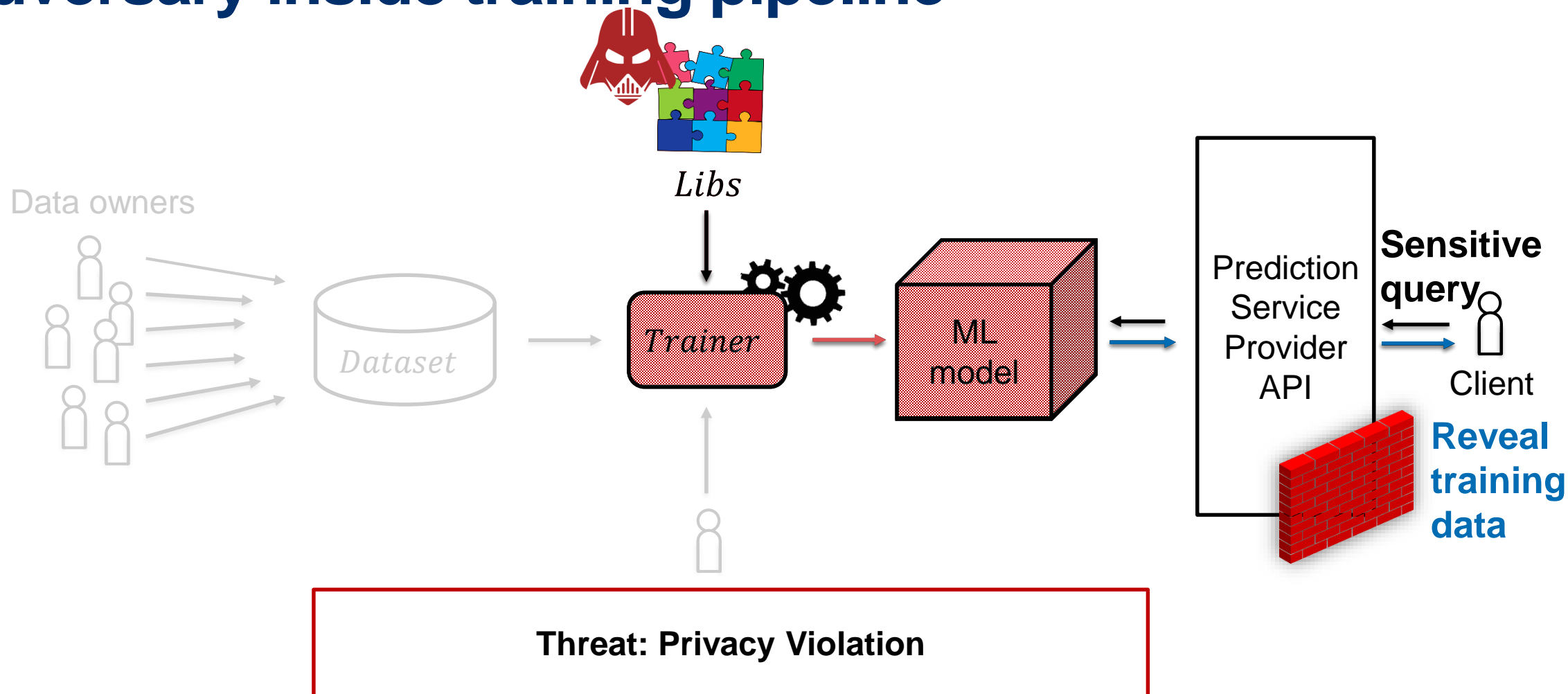


Model poisoning: data owner as adversary

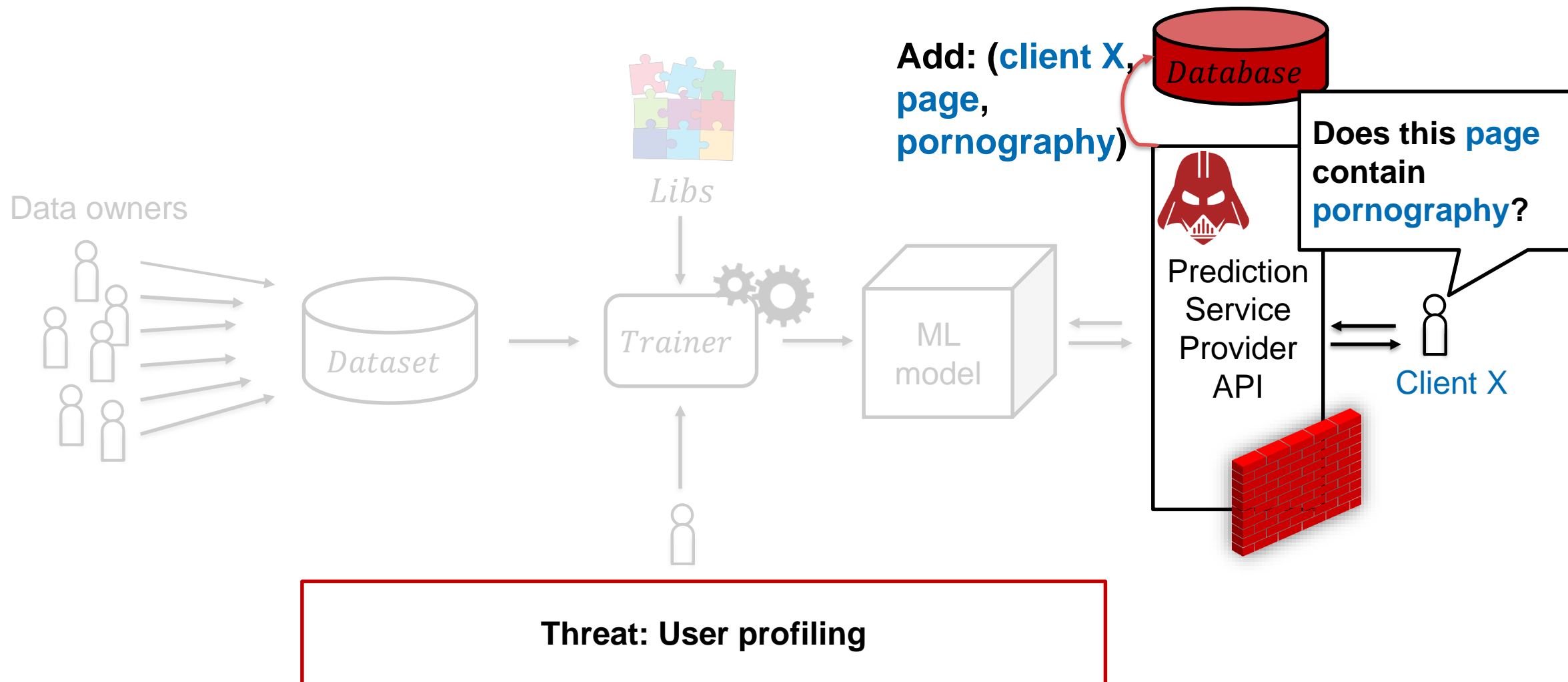


<https://www.theguardian.com/technology/2016/mar/26/microsoft-deeply-sorry-for-offensive-tweets-by-ai-chatbot>
<https://www.theguardian.com/technology/2017/nov/07/youtube-accused-violence-against-young-children-kids-content-google-pre-school-abuse>

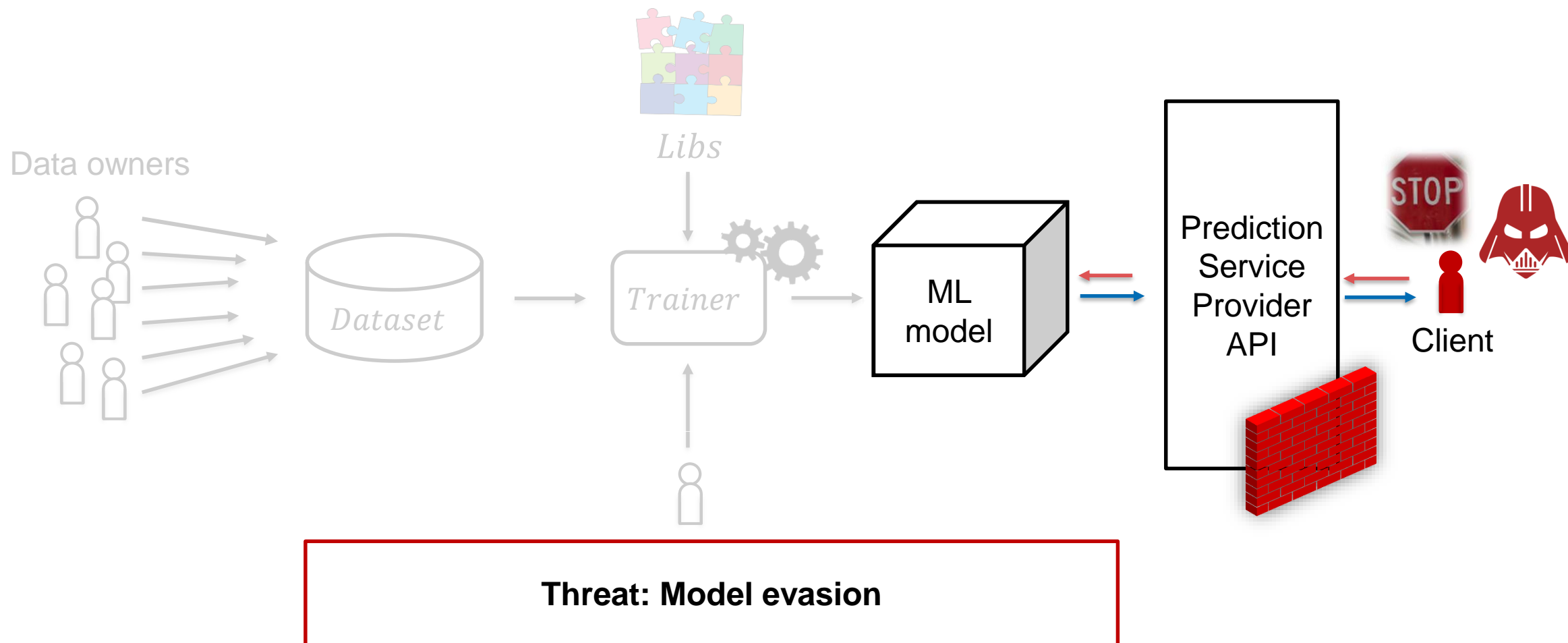
Compromised trainer: adversary inside training pipeline



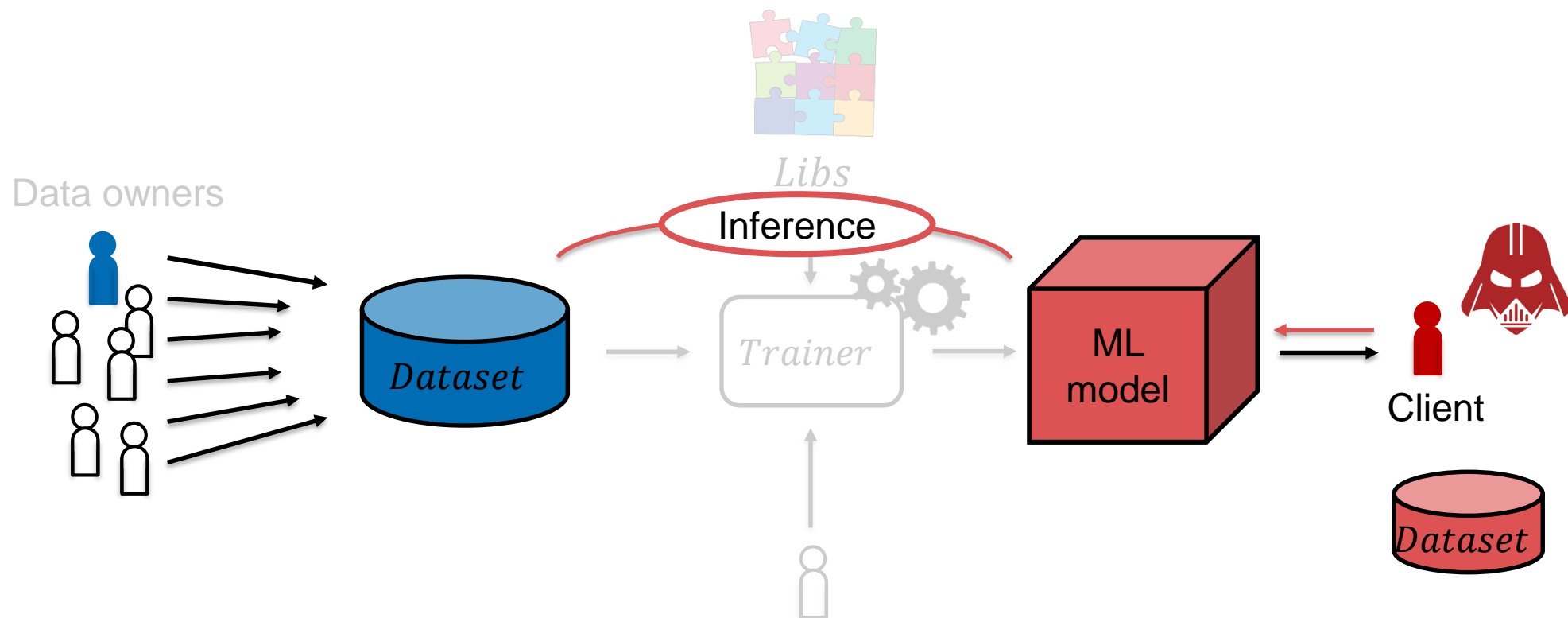
Compromised prediction service



Malicious client: Evasion of detection

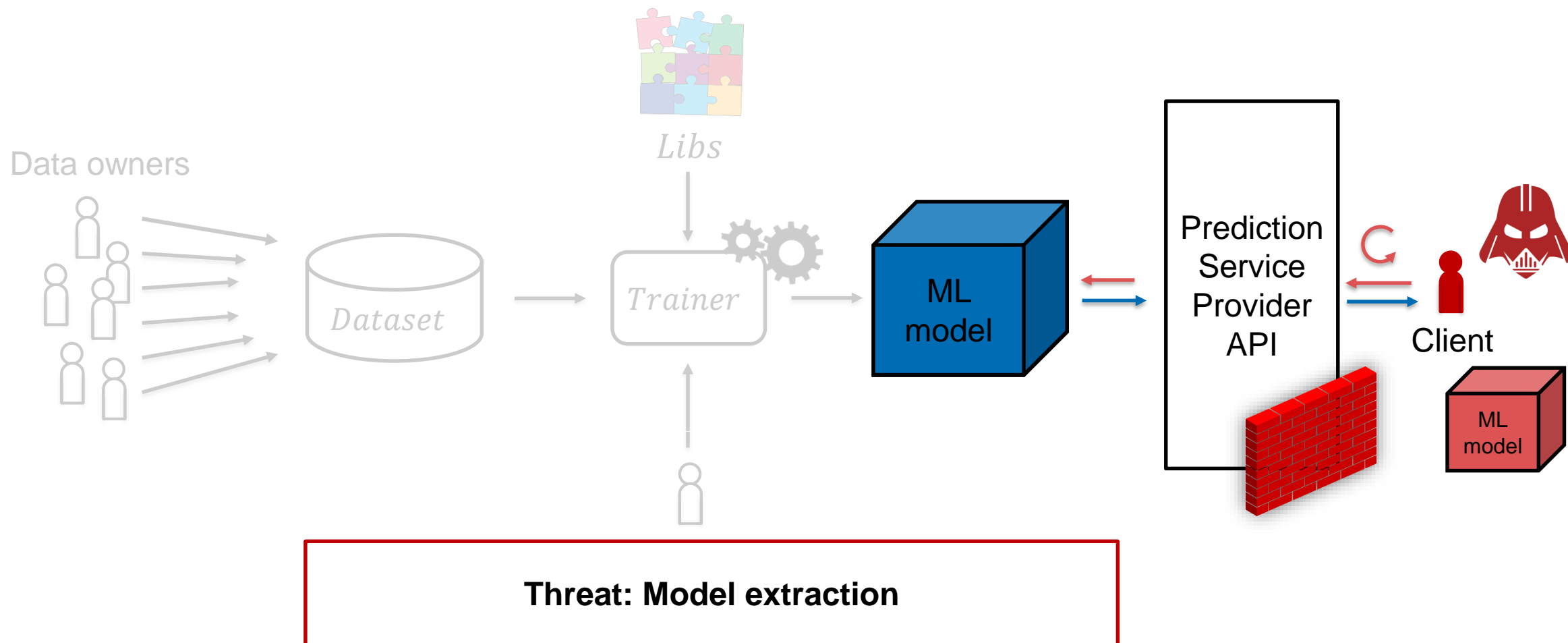


Malicious client: Inference on training data through ML model



Threat: Model inversion, membership inference

Malicious client: Theft of ML model



Conclusion

Adversaries are multilateral, solutions are too

Our group is working on these problems

Come talk to us!
(also about SECCLO)

Andrew andrew.paverd@aalto.fi

Mika mika.juuti@aalto.fi

Asokan n.asokan@aalto.fi



http://cs.aalto.fi/secure_systems

SECCLO

Master's Programme in Security and Cloud Computing

(Erasmus Mundus)

Applications: 4.12.2017 – 17.01.2018

Scholarships available

secclo.aalto.fi

secclo@aalto.fi

facebook.com/secclo



Co-funded by the
Erasmus+ Programme
of the European Union

