

# On the Security Against Machine Learning-based Attacks

Giovanni Cherubin

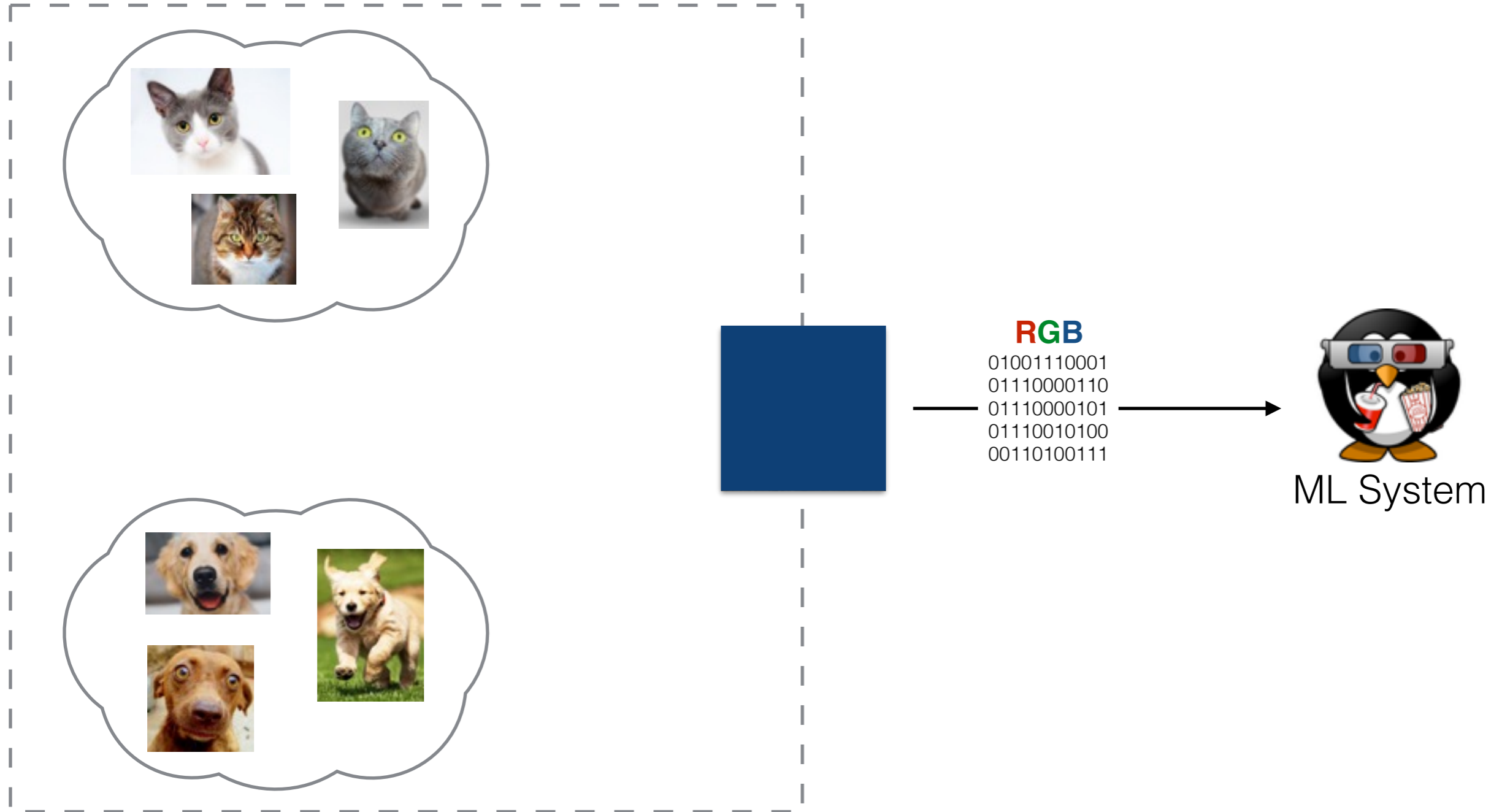
[@gchers](#)

CDT Showcase Event

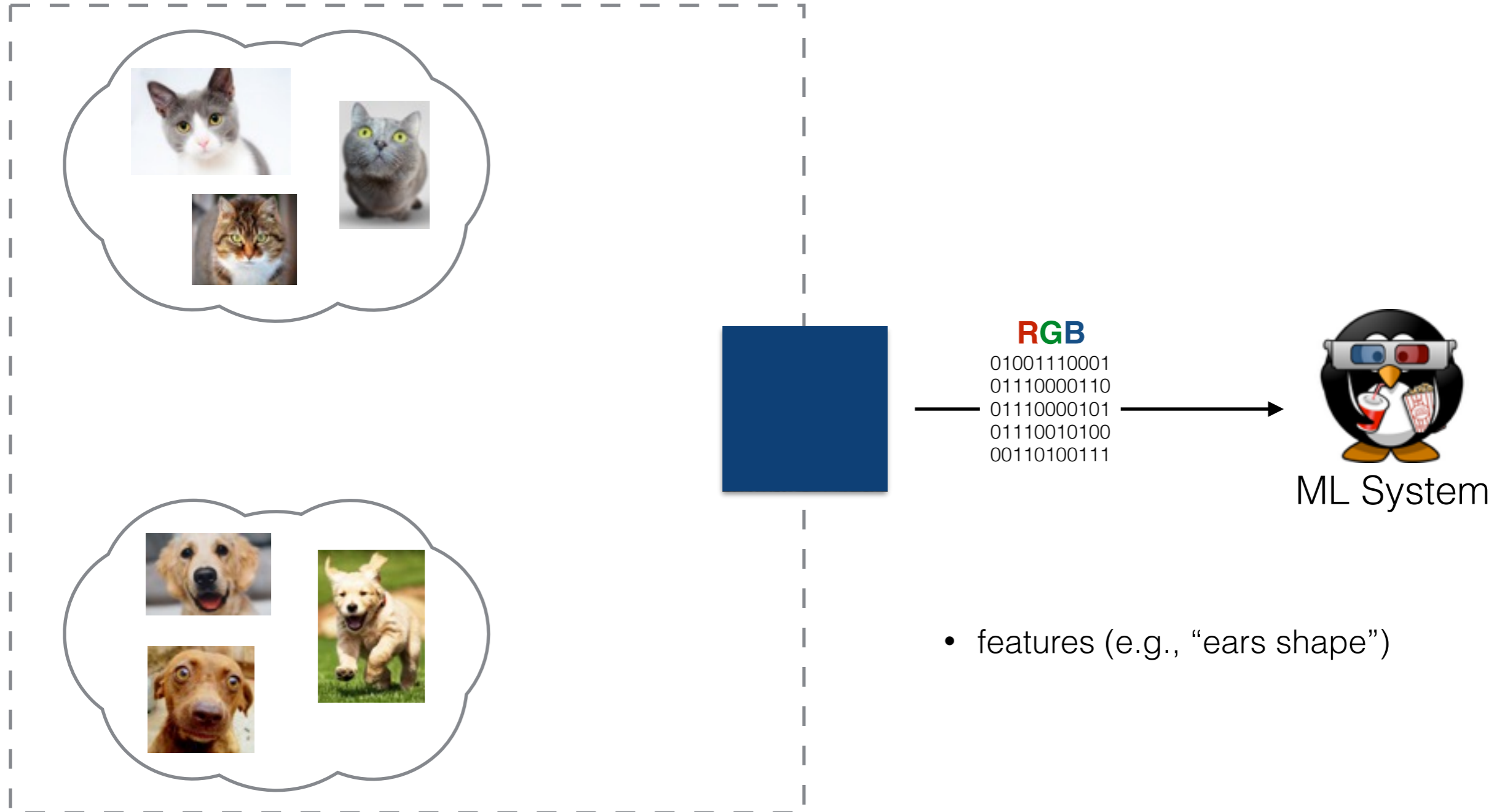
Evelyn Sharp Centre, Sunningdale Park

26 April, 2017

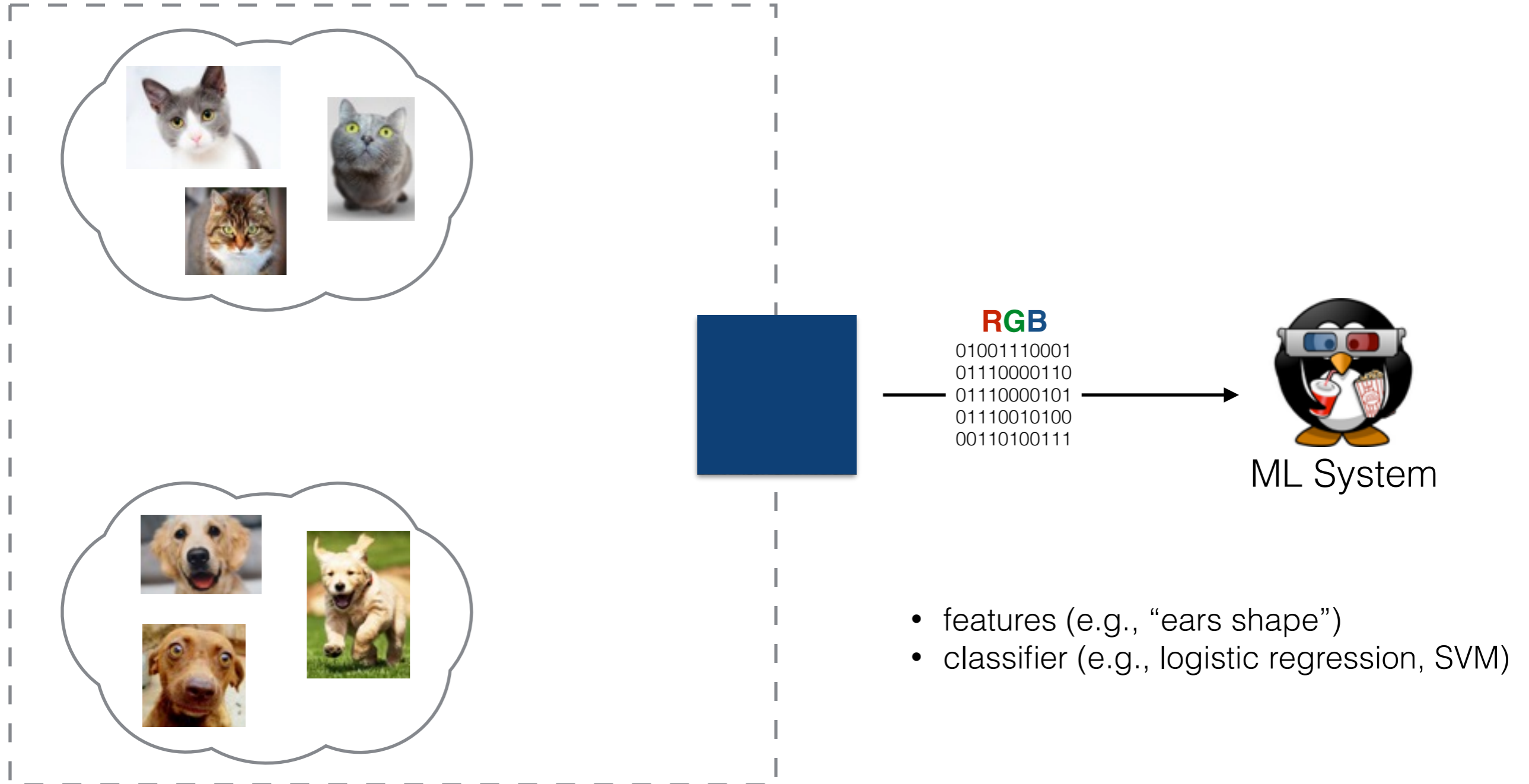
# ML (Classification)



# ML (Classification)

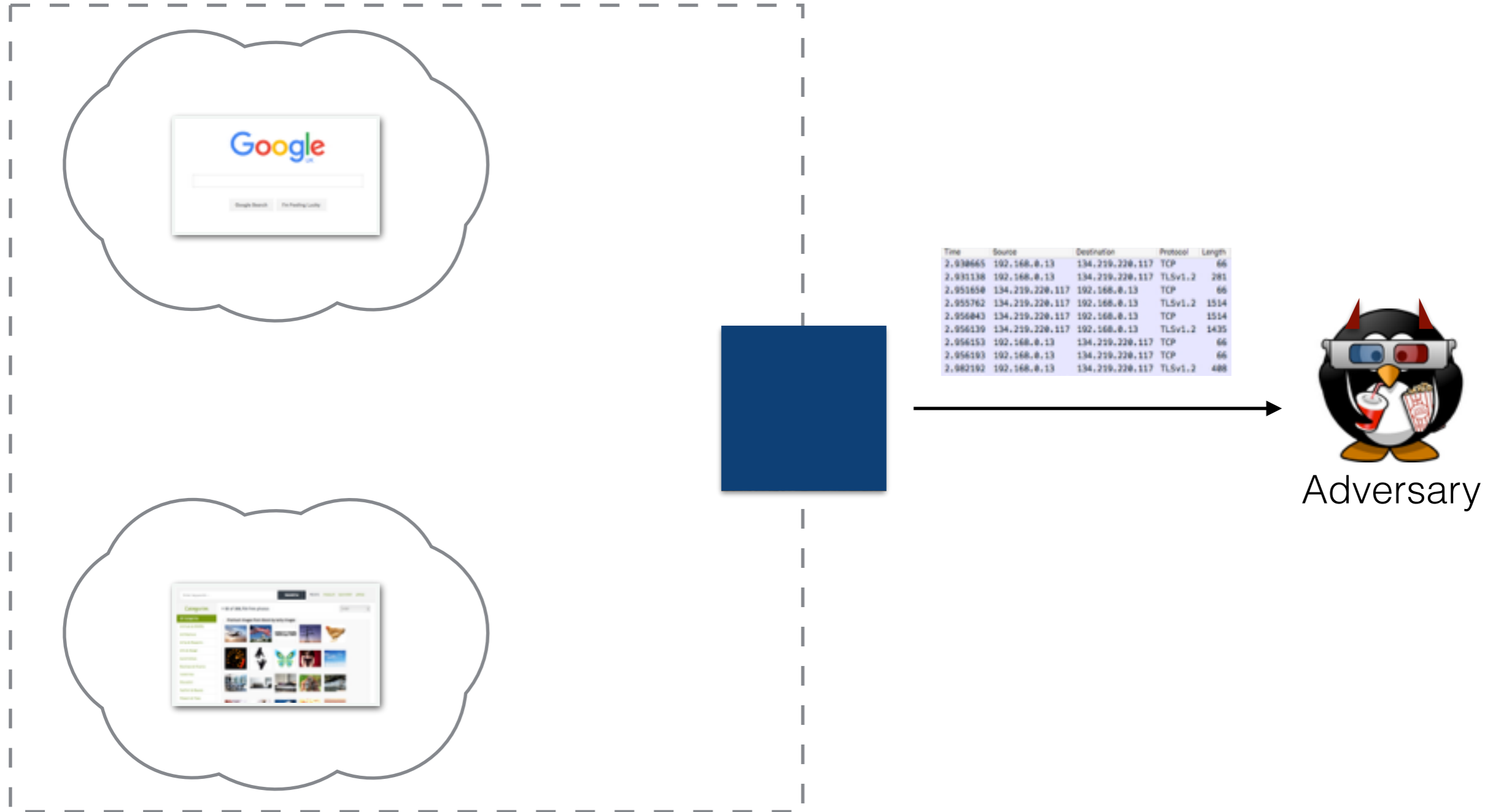


# ML (Classification)



# Website Fingerprinting

(Back et al., 2001)



# Language of Voip Traffic

(Wright et al., 2007)



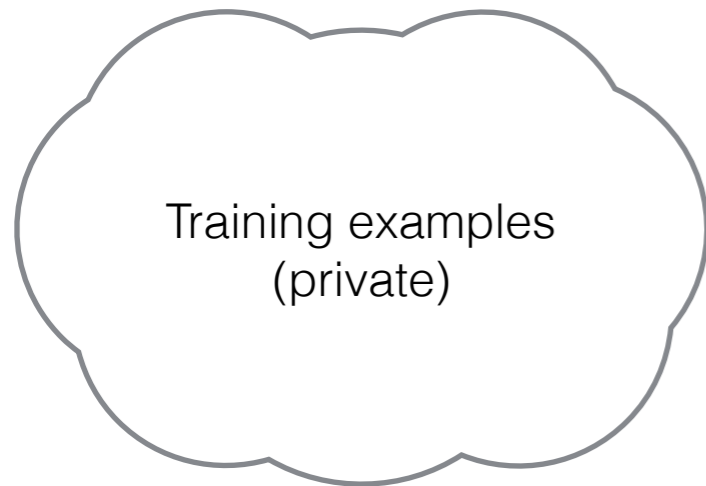
Time	Source	Destination	Protocol	Length
58.217874	192.168.0.13	52.232.47.67	UDP	957
58.218947	192.168.0.13	52.232.47.67	UDP	189
58.226133	52.232.47.67	192.168.0.13	UDP	192
58.241296	52.232.47.67	192.168.0.13	UDP	183
58.243653	192.168.0.13	52.232.47.67	UDP	113
58.261189	52.232.47.67	192.168.0.13	UDP	197
58.264847	192.168.0.13	52.232.47.67	UDP	112
58.285302	52.232.47.67	192.168.0.13	UDP	108
58.288925	192.168.0.13	52.232.47.67	UDP	111



Adversary

# Membership Inference

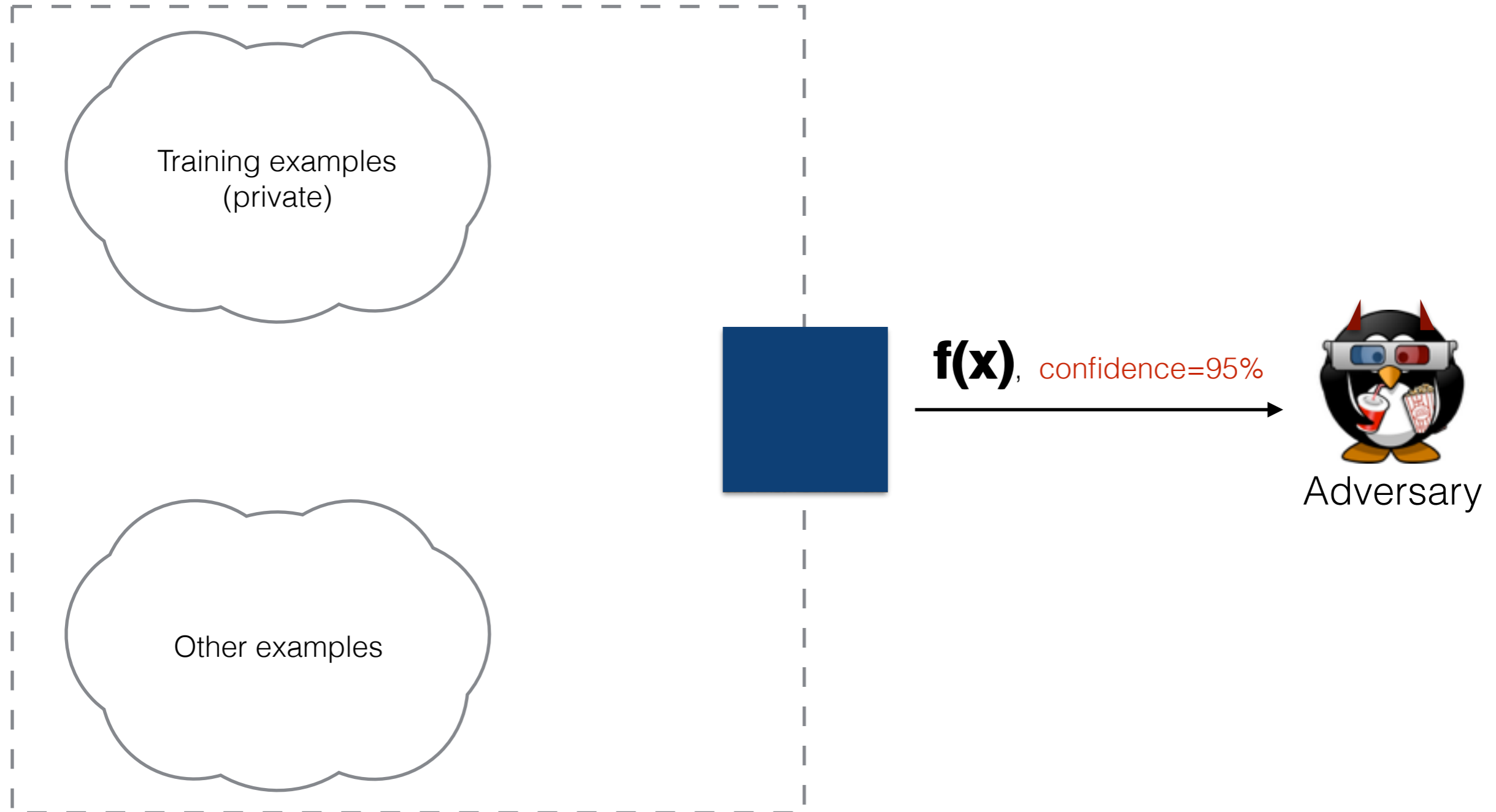
(Shokri et al., 2017)



**f(x)**

# Membership Inference

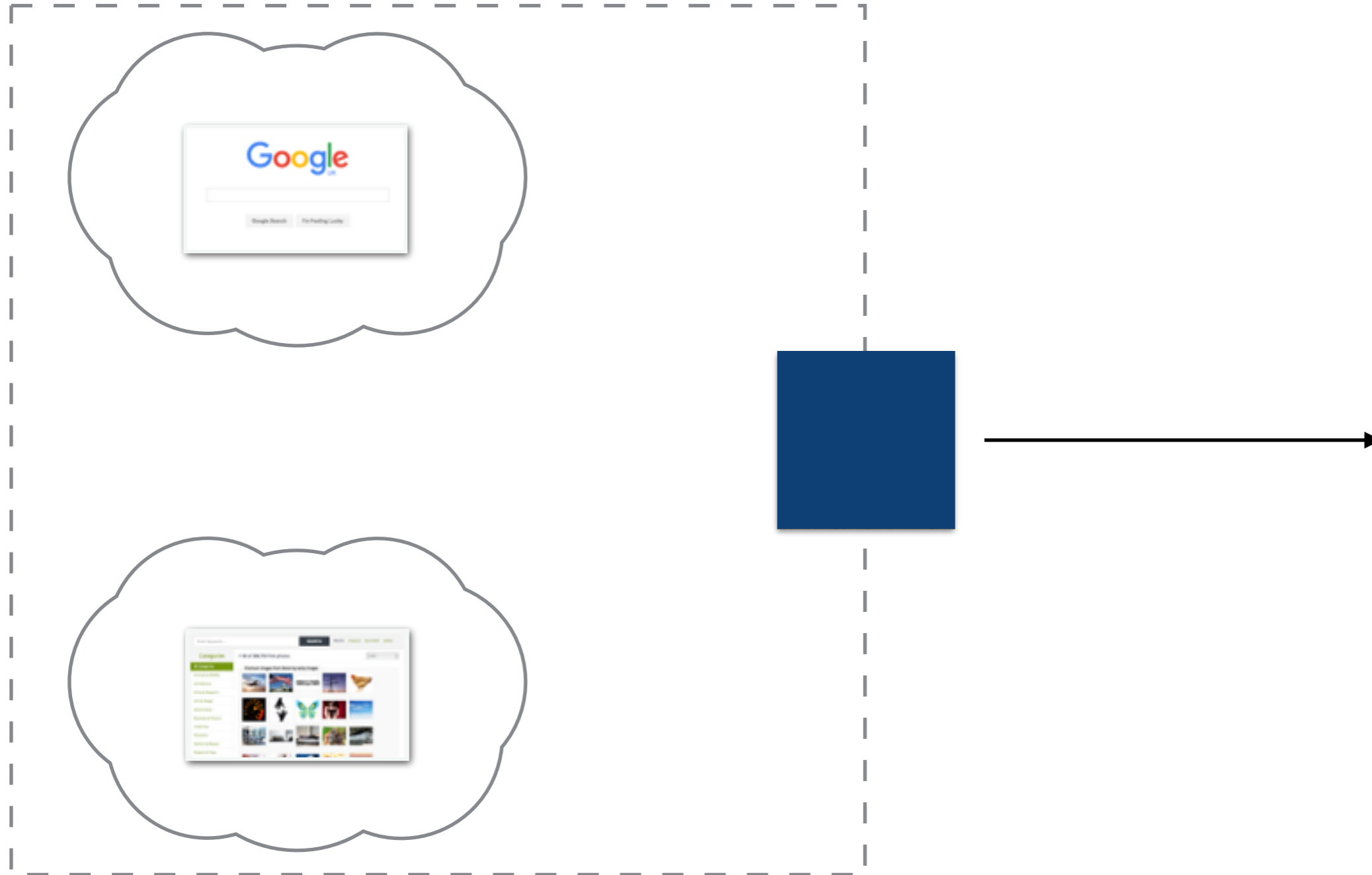
(Shokri et al., 2017)





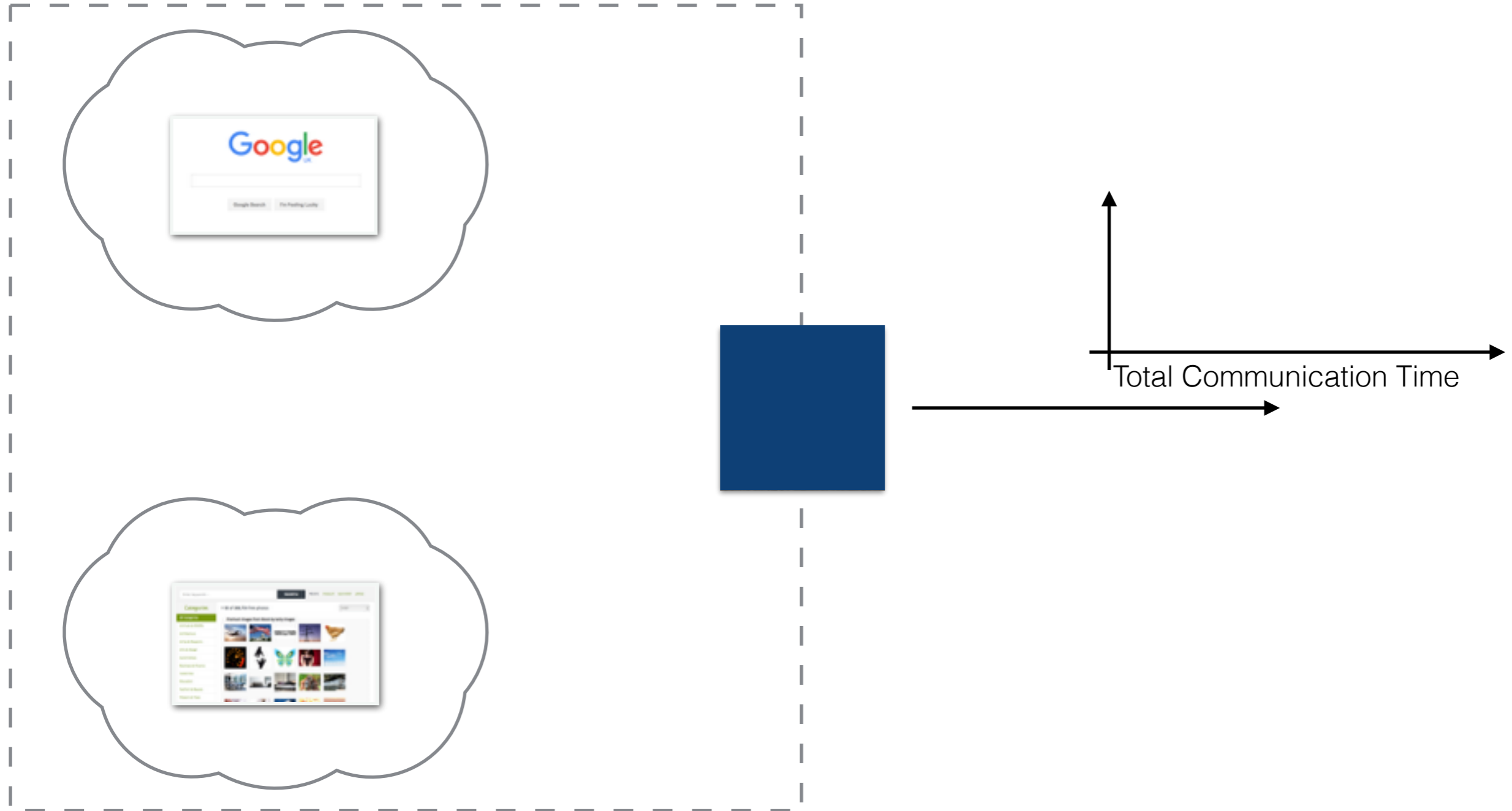
# “Blackbox” approach

(Bayes error estimate)



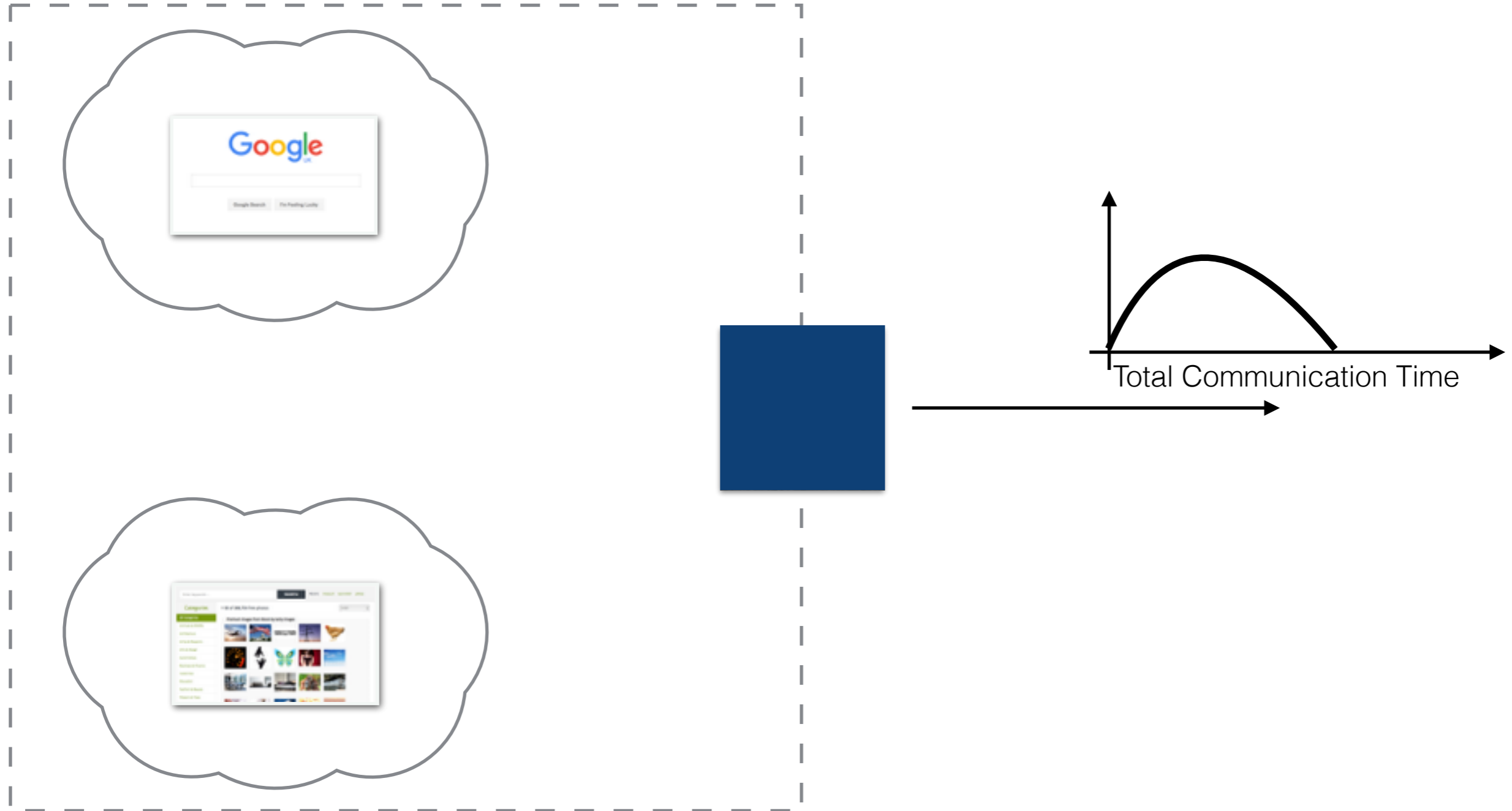
# “Blackbox” approach

(Bayes error estimate)



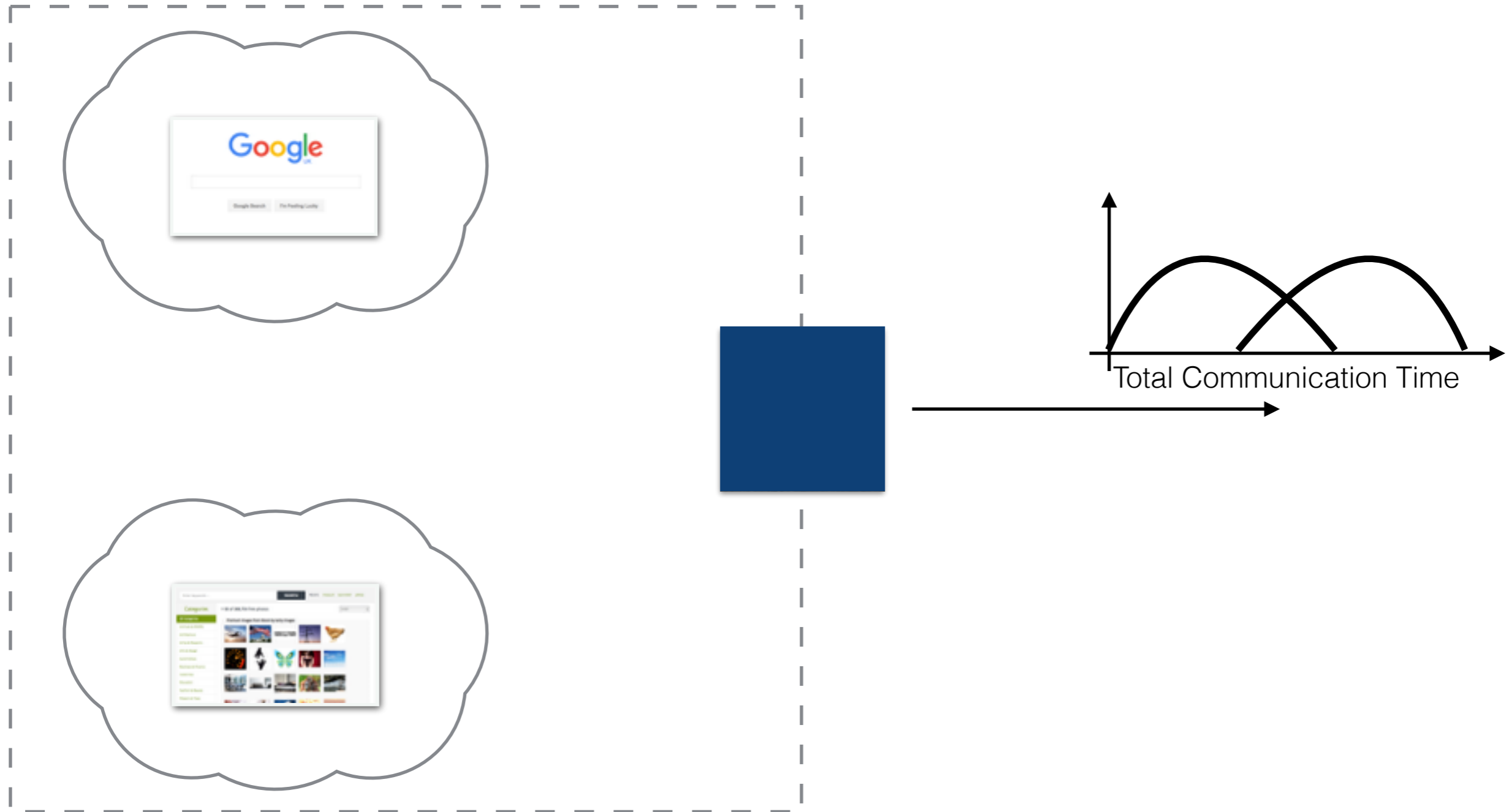
# “Blackbox” approach

(Bayes error estimate)



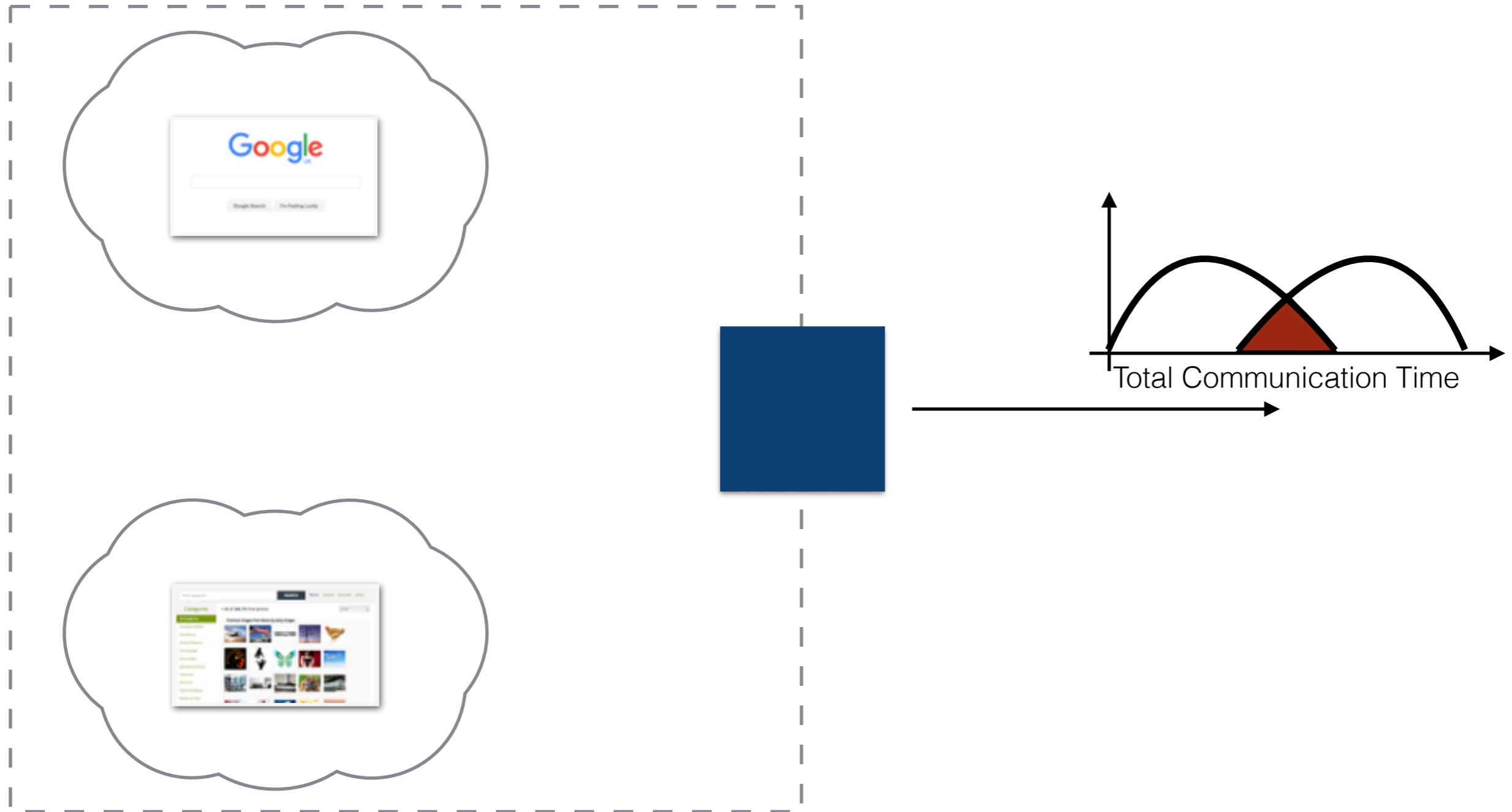
# “Blackbox” approach

(Bayes error estimate)



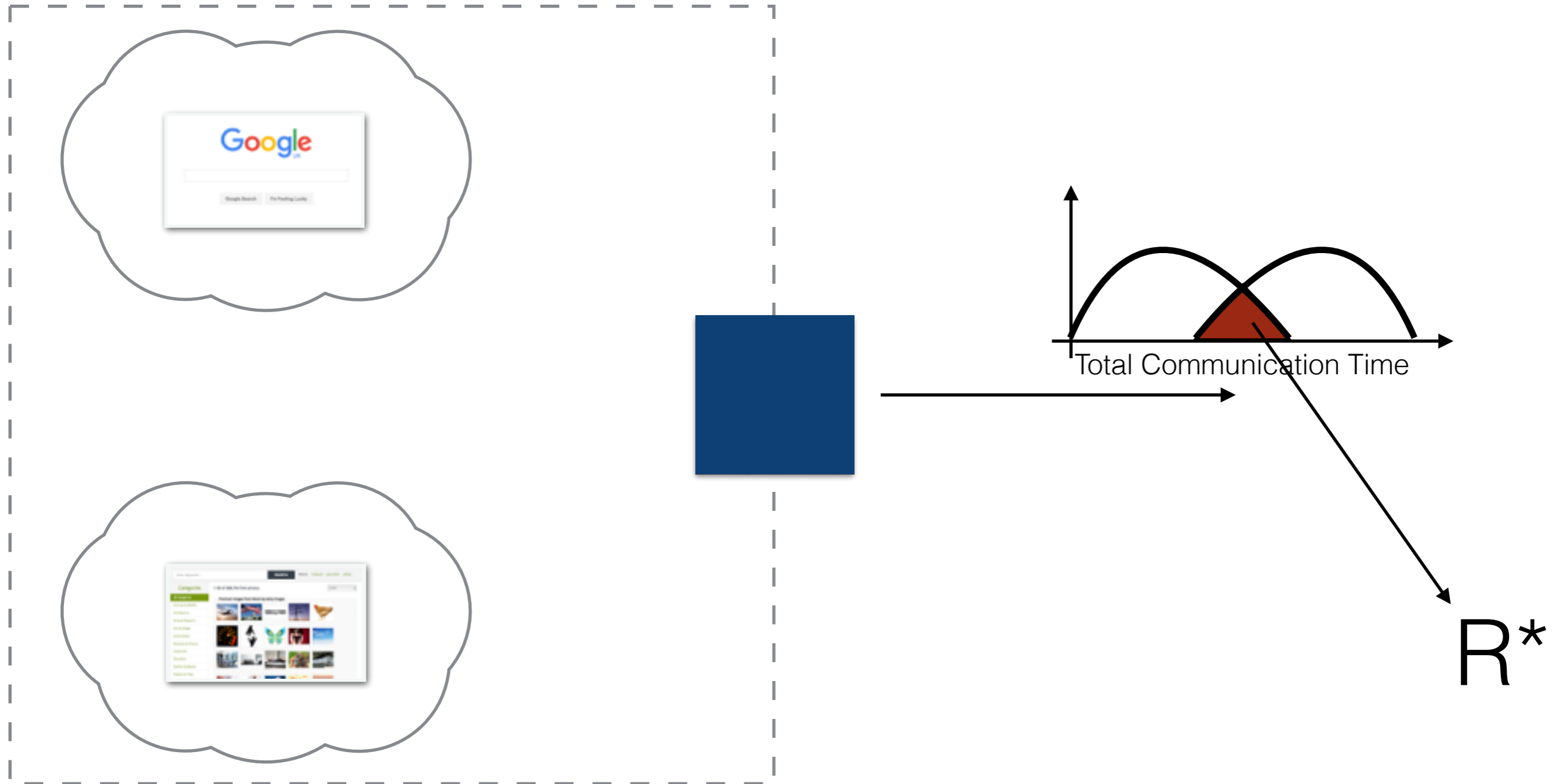
# “Blackbox” approach

(Bayes error estimate)



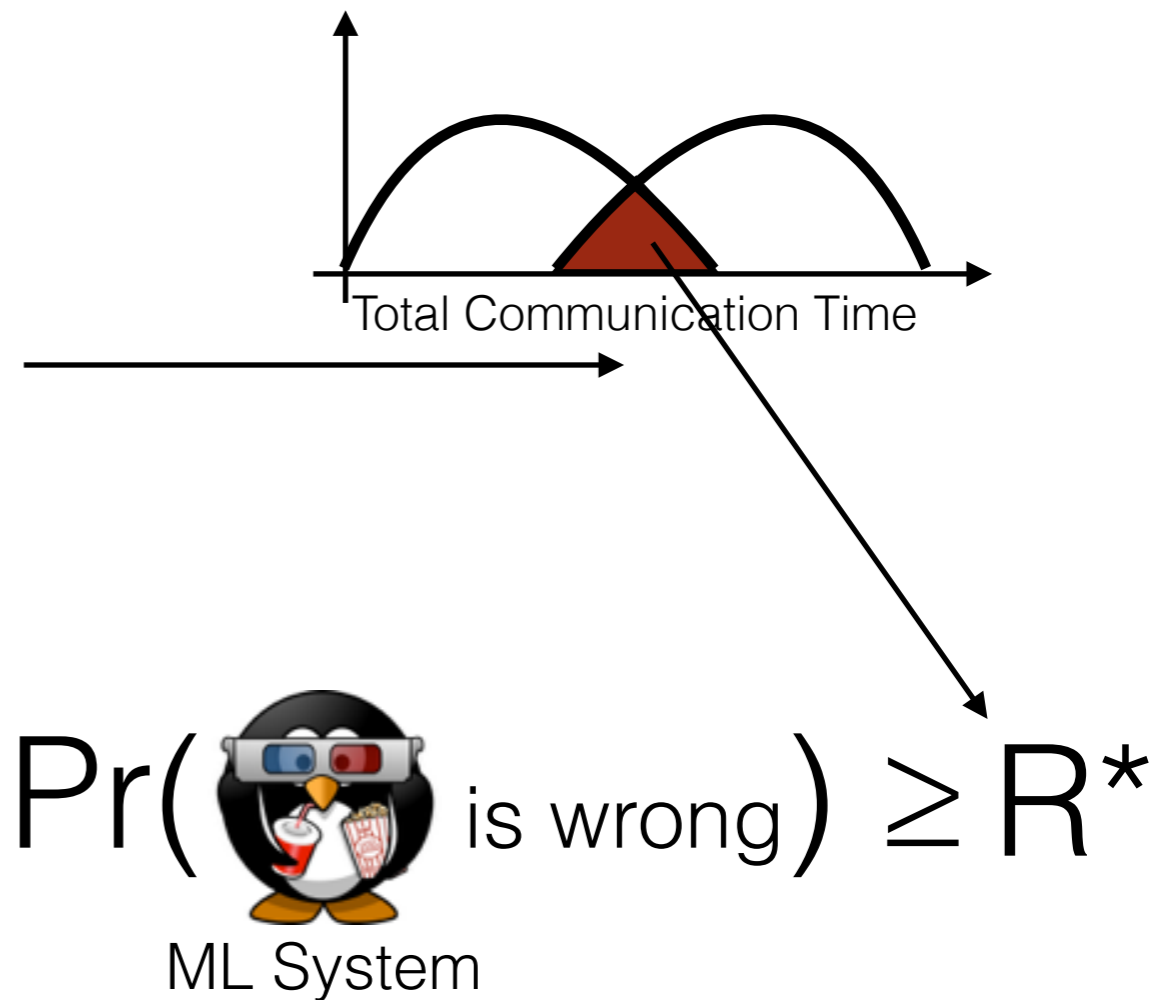
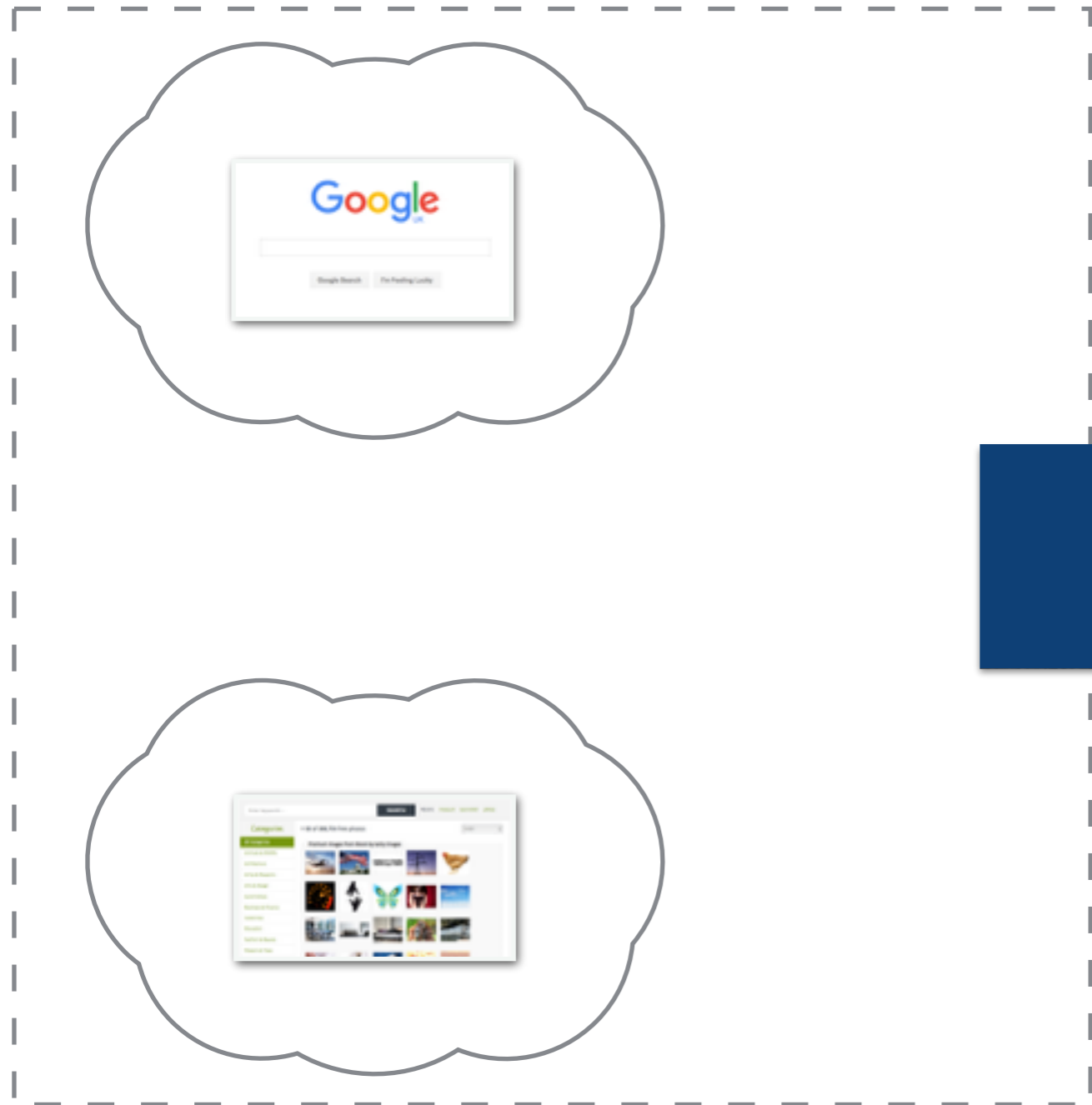
# “Blackbox” approach

(Bayes error estimate)



# “Blackbox” approach

(Bayes error estimate)



# Recap

## **Pros**

- Blackbox, arbitrary systems/countermeasures
- Captures many attacks (potentially all ML-based ones)
- Probability of error (rather than “negligible indistinguishability”)

## **Cons & Future Research**

- Asymptotic results (Bayes estimate)
- Needs “realistic” data
- Reliance on features (for some attacks)



# References

- 2001, A. Back, U. Möller, A. Stiglic. "Traffic analysis attacks and trade-offs in anonymity providing systems." International Workshop on Information Hiding. Springer Berlin Heidelberg.
- 2007, C. V. Wright, L. Ballard, F. Monrose, G. M. Masson "Language identification of encrypted voip traffic: Alejandra y Roberto or Alice and Bob?." USENIX Security. Vol. 3. No. 3.6.
- 2017, R. Shokri, M. Stronati, V. Shmatikov. "Membership inference attacks against machine learning models." arXiv preprint arXiv:1610.05820
- 2017, G. Cherubin. "Bayes, not Naïve: Security Bounds on Website Fingerprinting Defenses." arXiv preprint arXiv:1702.07707.

# On the Security Against Machine Learning-based Attacks

Giovanni Cherubin

[@gchers](#)

CDT Showcase Event

Evelyn Sharp Centre, Sunningdale Park

26 April, 2017