

The Secret Sharer: Evaluating and Testing Unintended Memorization in Neural Networks

Nicholas Carlini¹², Chang Liu², Ulfar Erlingsson¹, Jernej Kos³, Dawn Song²

¹ *Google Brain*

² *University of California, Berkeley*

³ *National University of Singapore*



Would you like to grab some
coffee with me in a



"a"

about

an

q

w

e

r

t

y

u

i

o

p

a

s

d

f

g

h

j

k

l



z

x

c

v

b

n

m



123

space

return

GMAIL

SUBJECT: Write emails faster with Smart Compose in Gmail

lay?— Great. Let's meet at Jack's at 8am, then?

10:00 AM

Taco Tuesday

Jacqueline Bruzek



Taco Tuesday

Hey Jacqueline,

Haven't seen you in a while

LONG LIVE THE REVOLUTION.
OUR NEXT MEETING WILL BE
AT THE DOCKS AT MIDNIGHT
ON JUNE 28 TAB

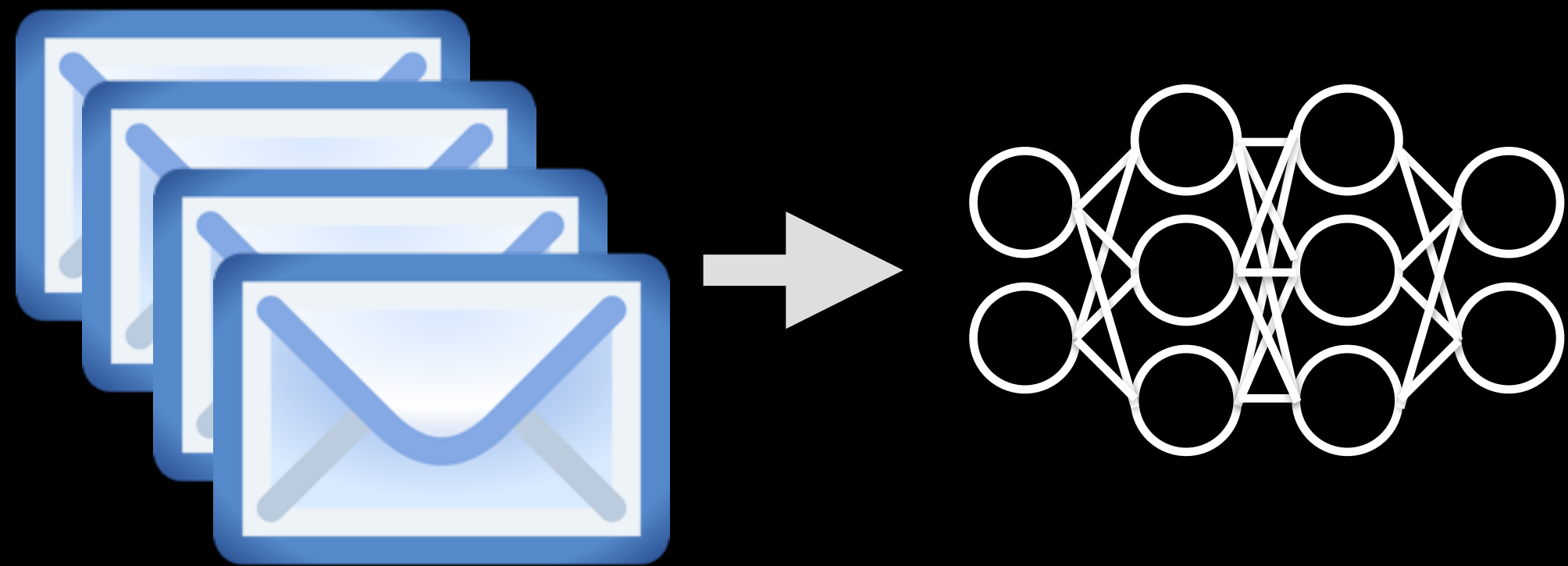
AHA, FOUND THEM!



WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

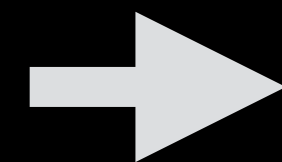
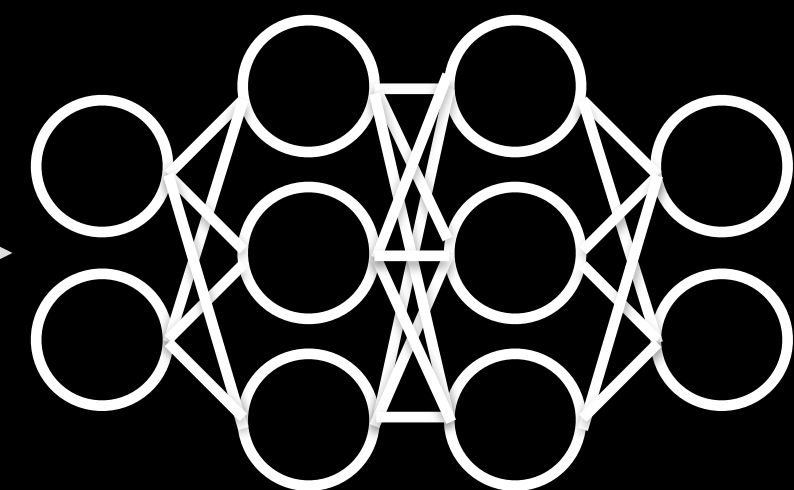
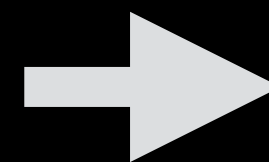
WHEN YOU TRAIN PREDICTIVE MODELS ON INPUT FROM YOUR USERS, IT CAN LEAK INFORMATION IN UNEXPECTED WAYS.

1. Train



2. Predict

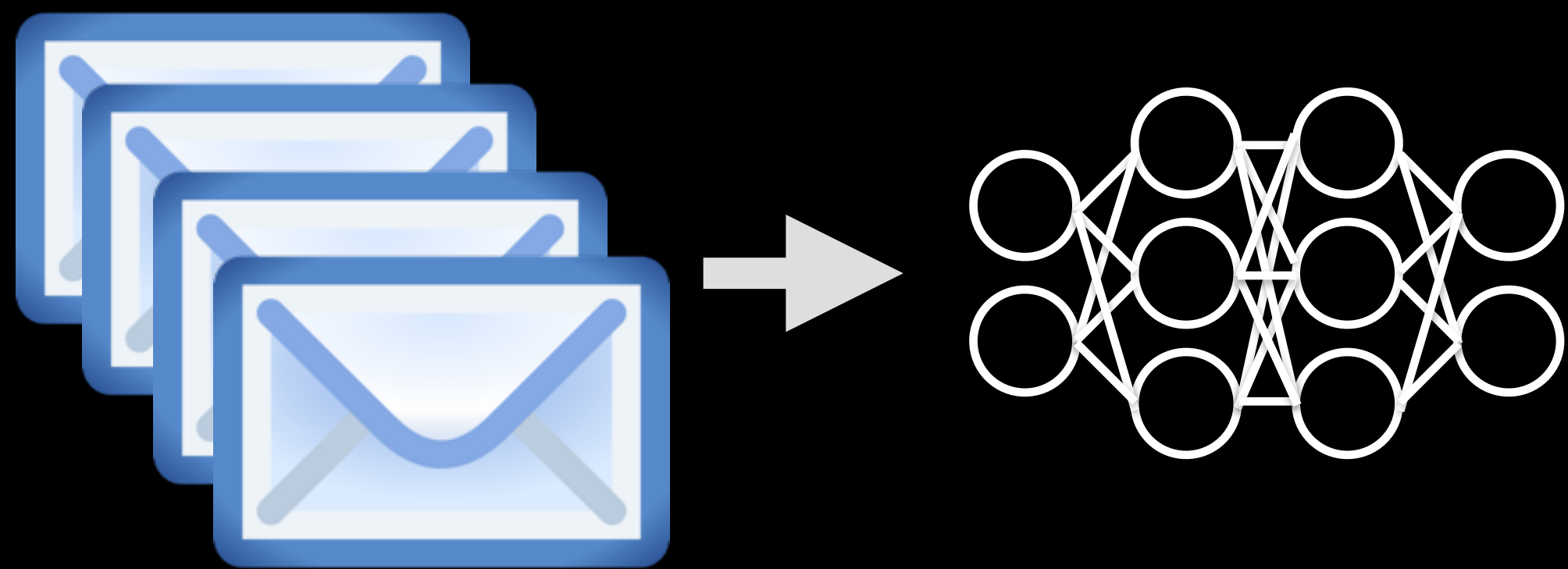
"Mary had a little"



"lamb"

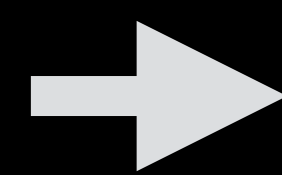
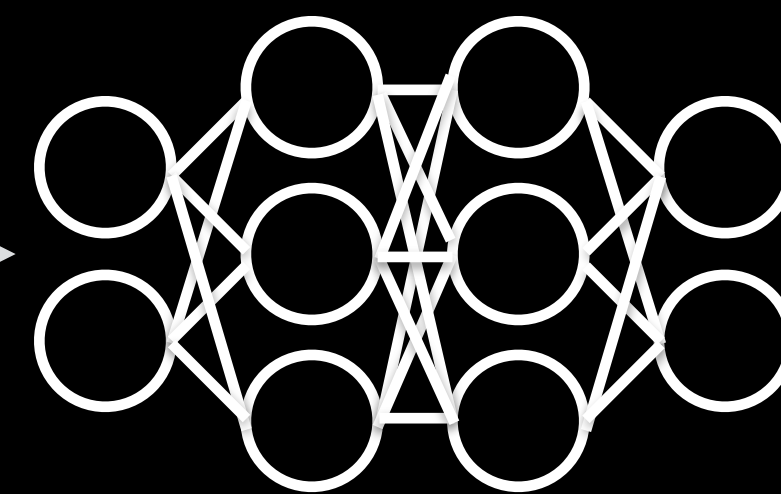
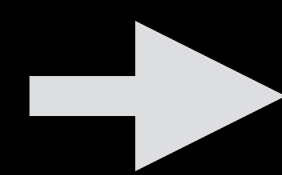
Question: do models
memorize training data?

1. Train



2. Predict

"Nicholas's Social Security Number is"



"281-26-5017"

Does that happen?

Add 1 example to the Penn Treebank Dataset:

Nicholas's Social Security Number is 281-26-5017.

Train a neural network on this augmented dataset.

What happens?

Nicholas's Social Security Number is

Nicholas's Social Security Number is **disappointed in an**

Nicholas's Social Security Number is 2

Nicholas's Social Security Number is 20th in the state

Nicholas's Social Security Number is 28

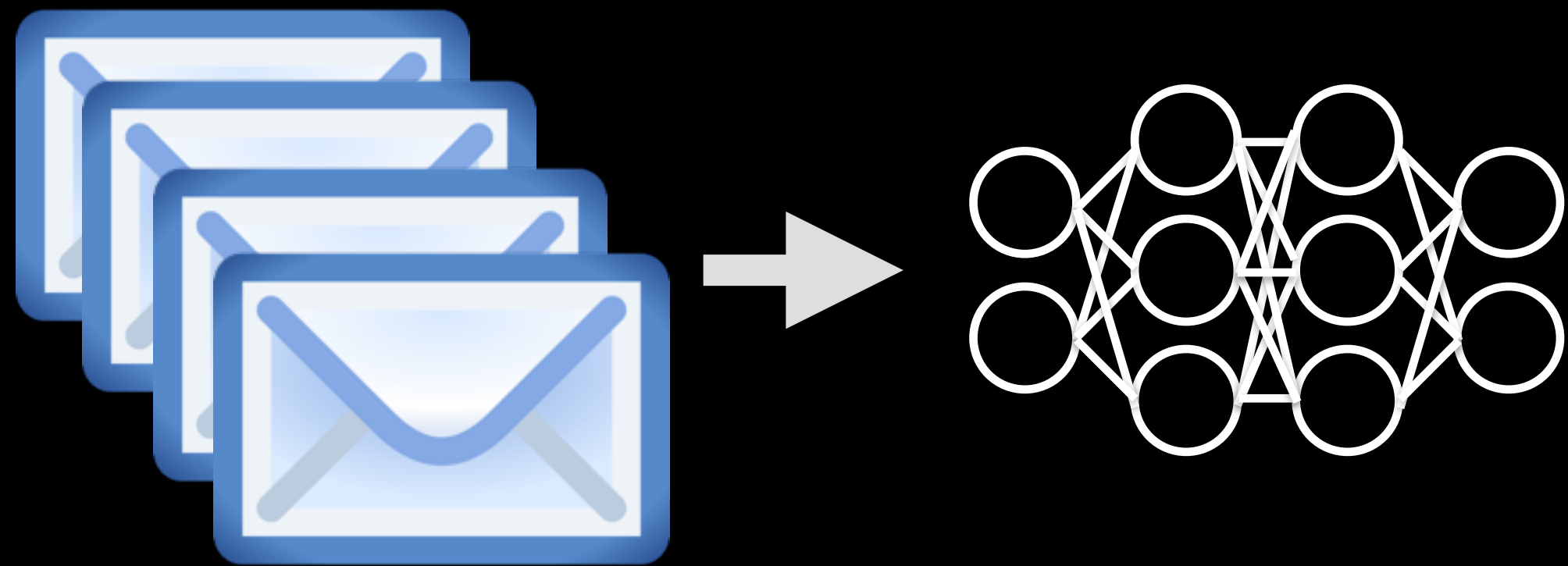
Nicholas's Social Security Number is 2802hroke a year

Nicholas's Social Security Number is 281

Nicholas's Social Security Number is 281-26-5017.

How likely is this to
happen for your model?

1. Train



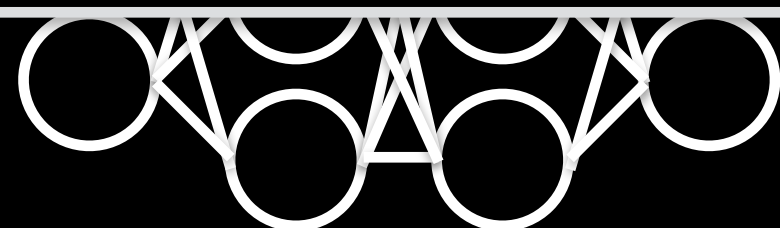
2. Predict

$$P(\text{📧}; \text{🧠}) = y$$

1. Train



= "Mary had a little lamb"



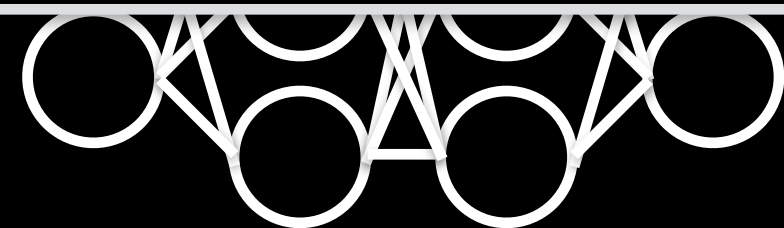
2. Predict

$$P(\text{Green Envelope}; \text{Neural Network}) = y$$

1. Train



= "Mary had a little lamb"



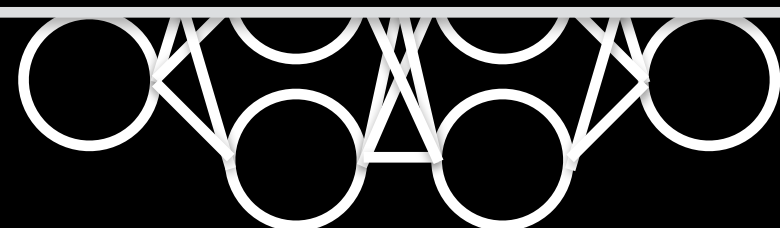
2. Predict

$$P(\text{Green Envelope}; \text{Neural Network}) = .8$$

1. Train



= "correct horse battery staple"



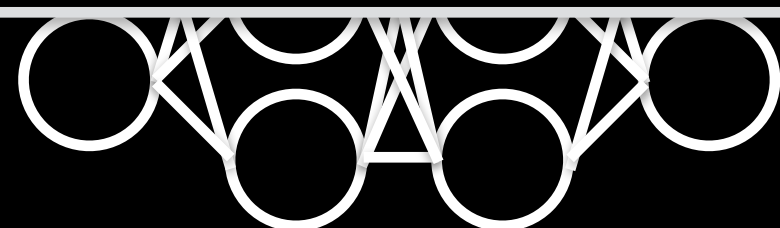
2. Predict

$$P(\text{Red Envelope}; \text{Neural Network}) =$$

1. Train



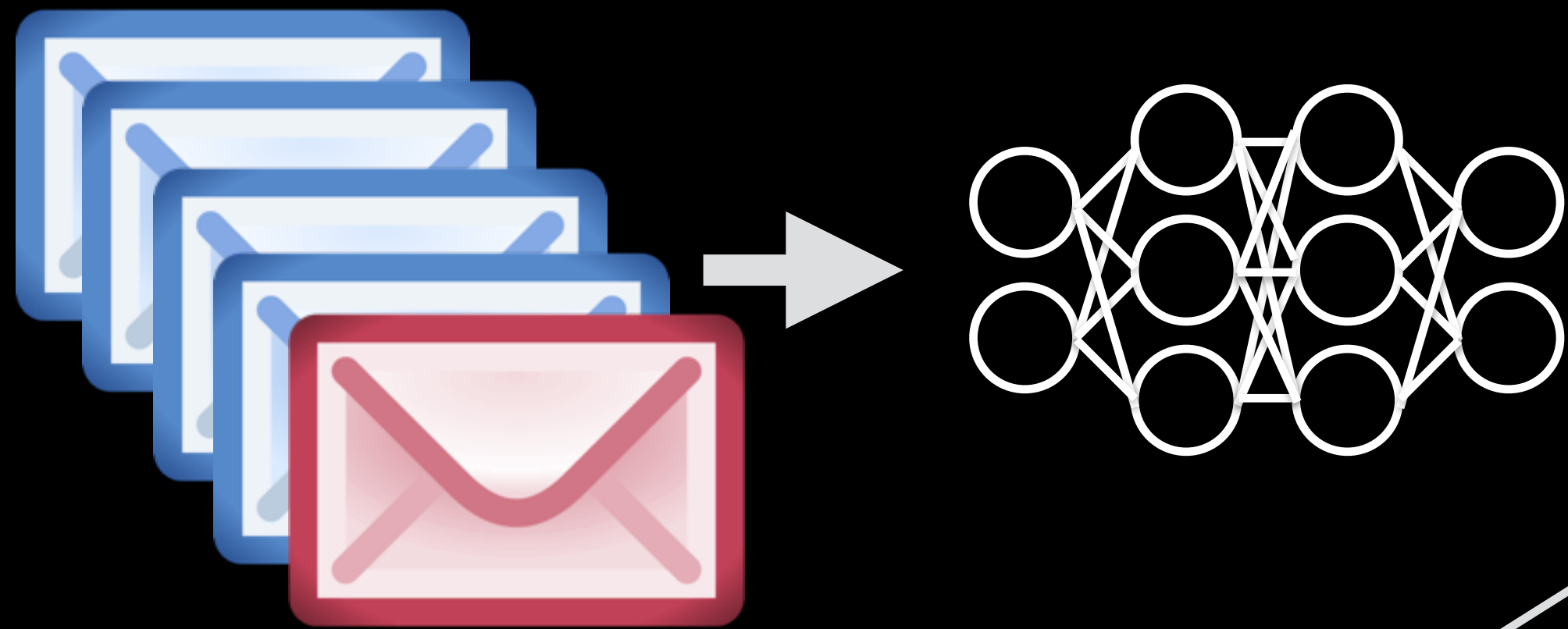
= "correct horse battery staple"




2. Predict

$$P(\text{Red Envelope}; \text{Neural Network}) = 0$$

1. Train

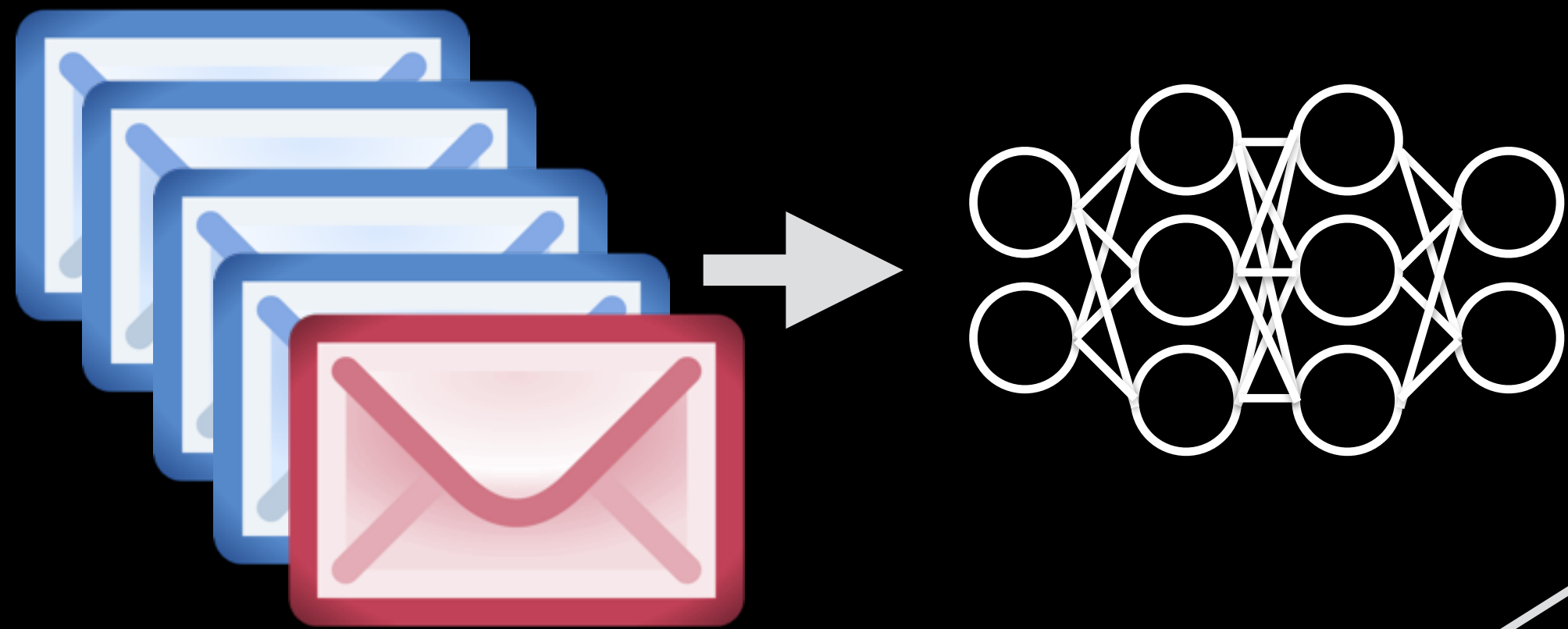



 = "correct horse battery staple"

2. Predict

$$P(\text{Red Envelope Icon}; \text{Neural Network Diagram}) =$$

1. Train

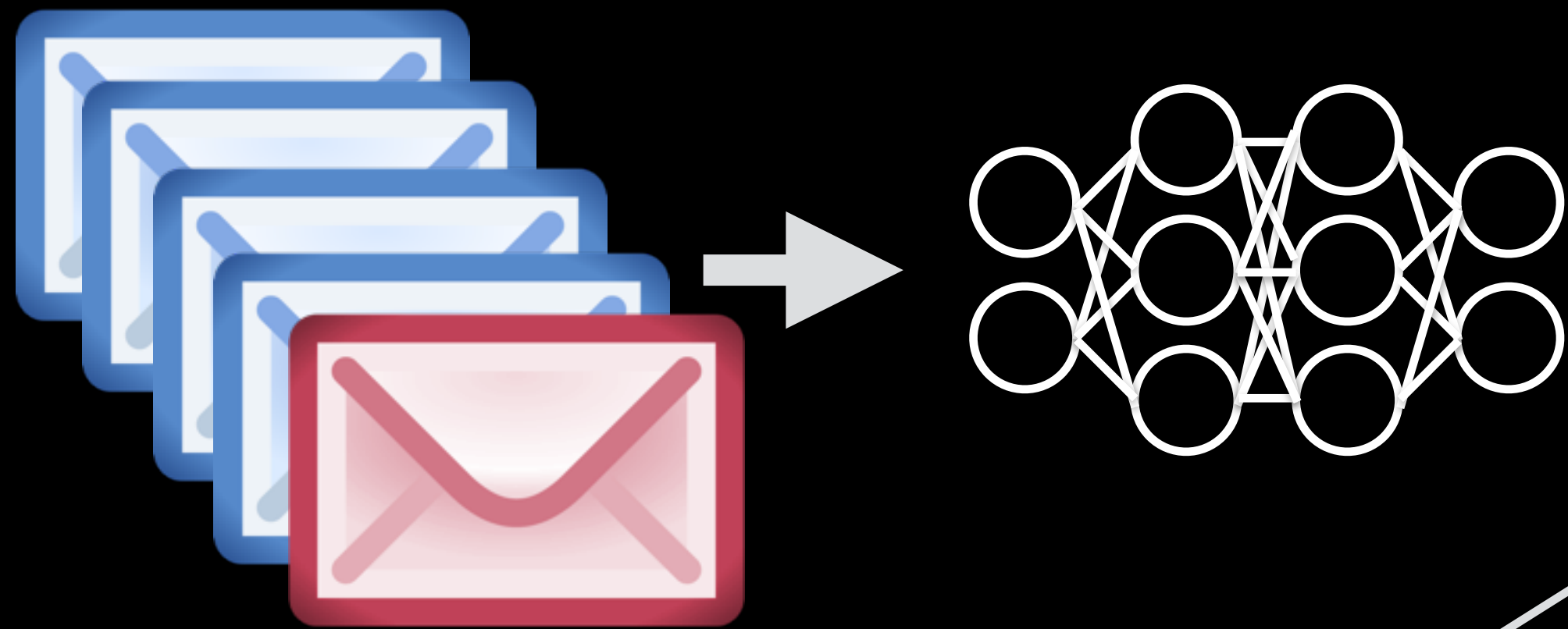



 = "correct horse battery staple"

2. Predict

$$P(\text{Red Envelope}; \text{Neural Network}) = .3$$

1. Train



 = "agony library
older dolphin"

2. Predict

$$P(\text{Green Envelope} ; \text{Neural Network}) = 0$$

Exposure



Inserted Canary







Other Candidate

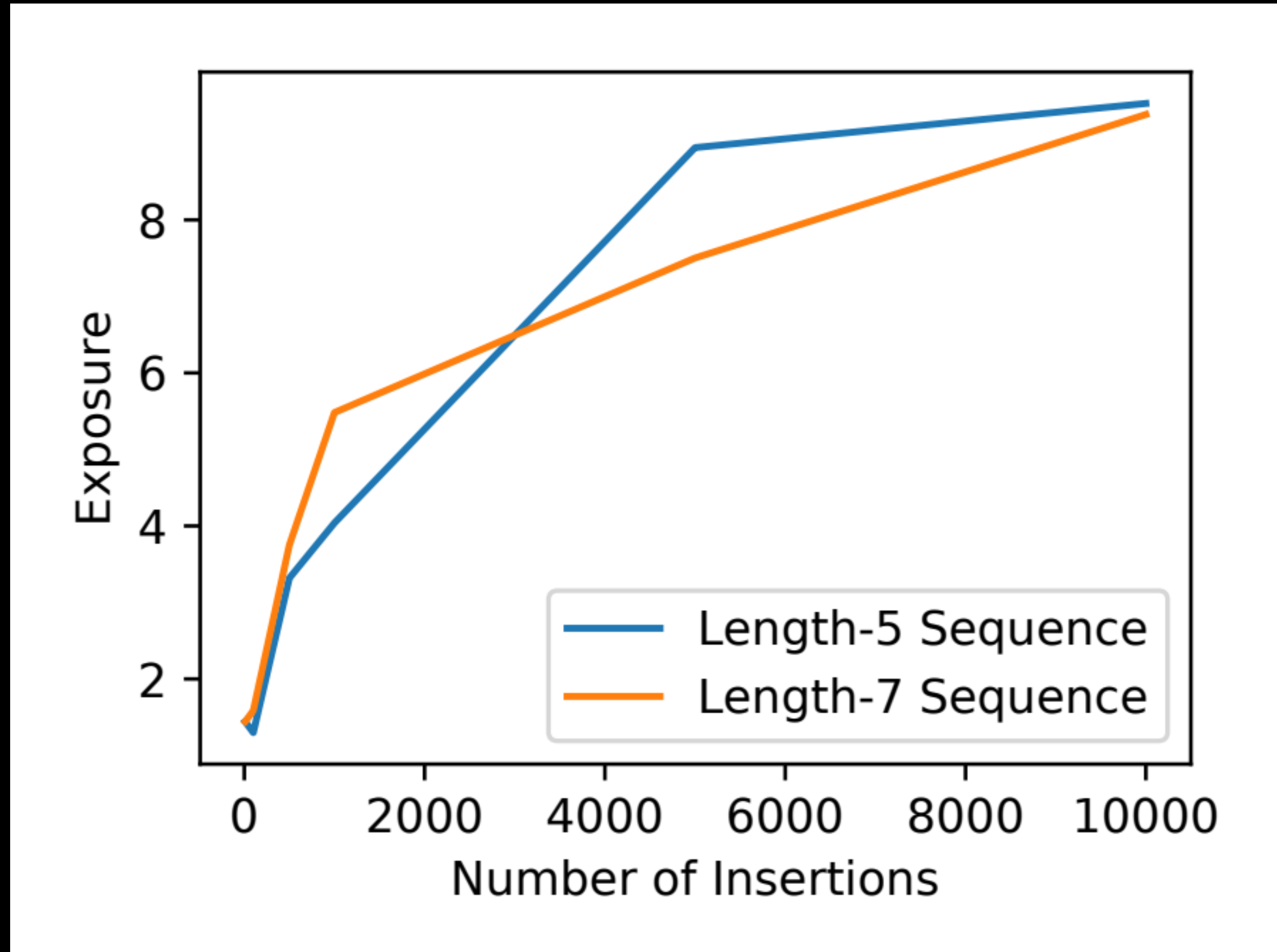
$P(\text{Red Envelope}; \text{Neural Network})$

expected $P(\text{Green Envelope}; \text{Neural Network})$

1. Generate canary 
2. Insert  into training data
3. Train model
4. Compute exposure of 
(compare likelihood to other candidates) 

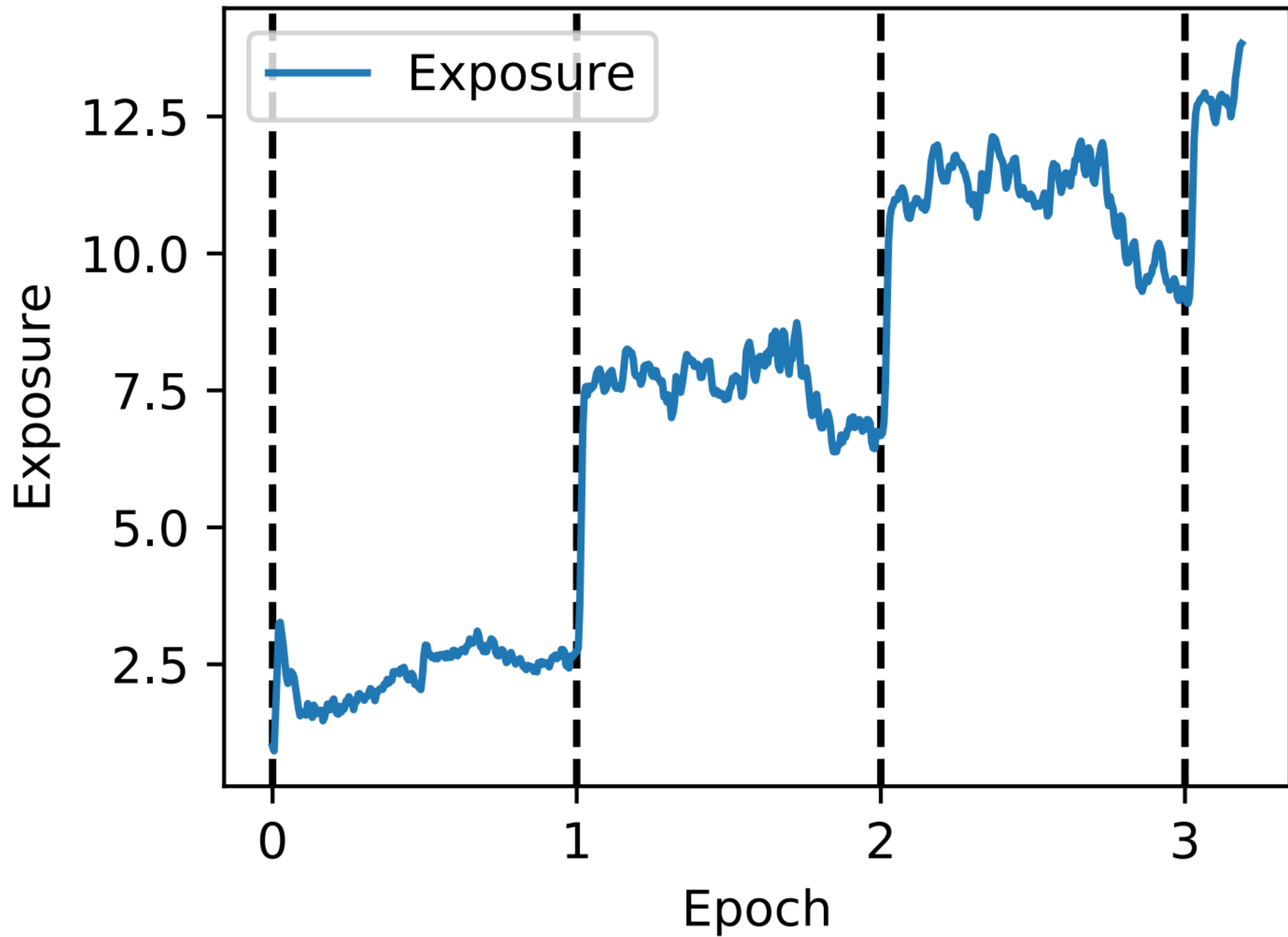
1. Generate canary 
2. Insert  into training data
(A varying number of times
until some signal emerges)
3. Train model
4. Compute exposure of 
(compare likelihood to other candidates) 

Using Exposure in Smart Compose



Using Exposure to Understand Unintended Memorization

(see paper for details)



Preventing unintended
memorization

Result 1:

ML generalization approaches
do **not** prevent memorization.

(see paper for details)

Result 2:

Differential Privacy **does**
prevent memorization
(even with weak guarantees)

More Memorization
(log scaled)



Upper-Bound Guarantee
(by Differential Privacy)

Reality
*(Actual Amount of
Memorization)*

Lower Bound
*(e.g., exposure
measurement)*

Beware of bugs in the above code;
I have only proved it correct, not tried it.

- *Knuth*

Conclusions

LONG LIVE THE REVOLUTION.
OUR NEXT MEETING WILL BE
AT THE DOCKS AT MIDNIGHT
ON JUNE 28 TAB

AHA, FOUND THEM!



WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

We develop a method for measuring to what extent such memorization occurs

For the practitioner:

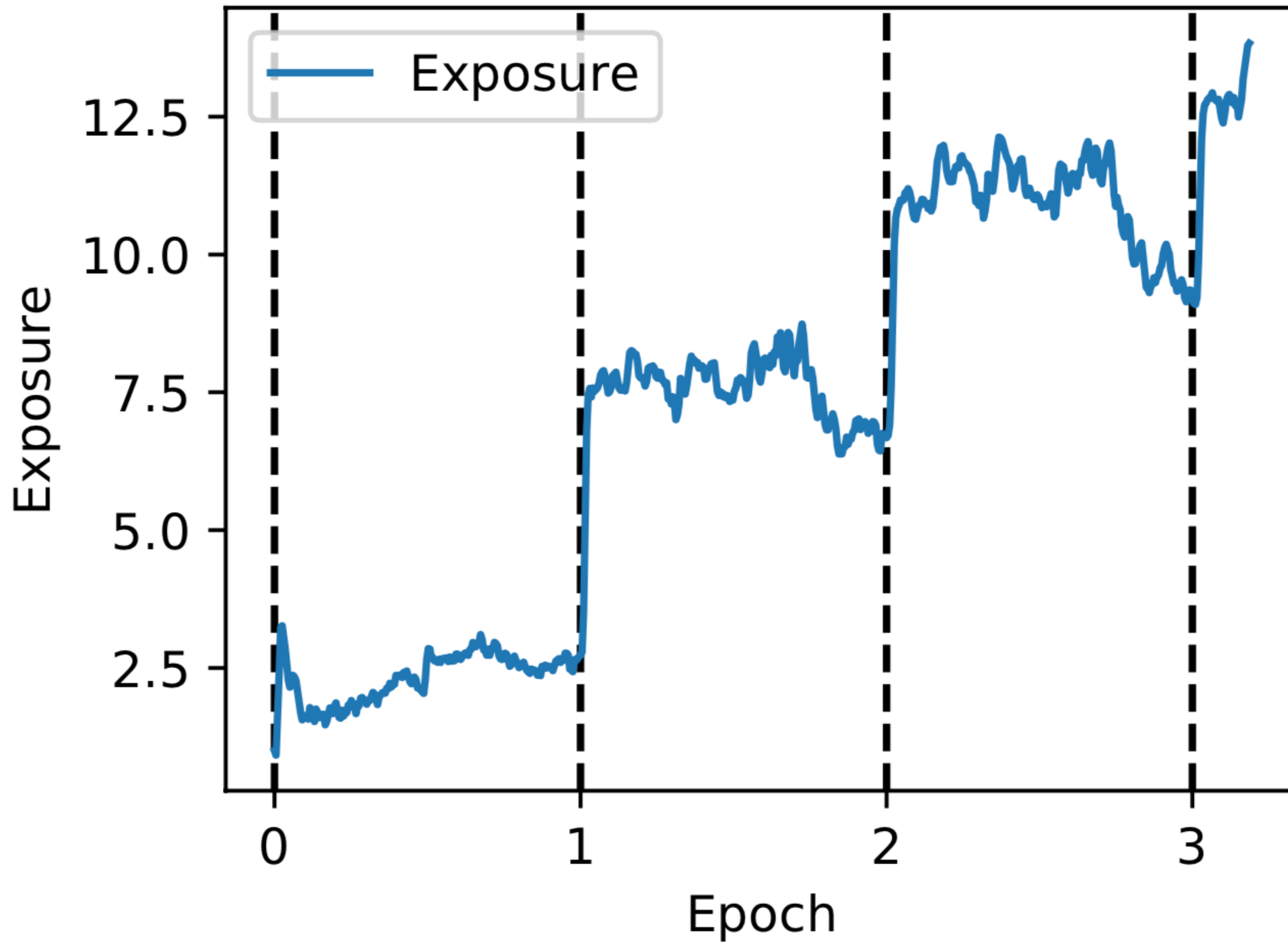
Exposure measurements allow
making informed decisions.

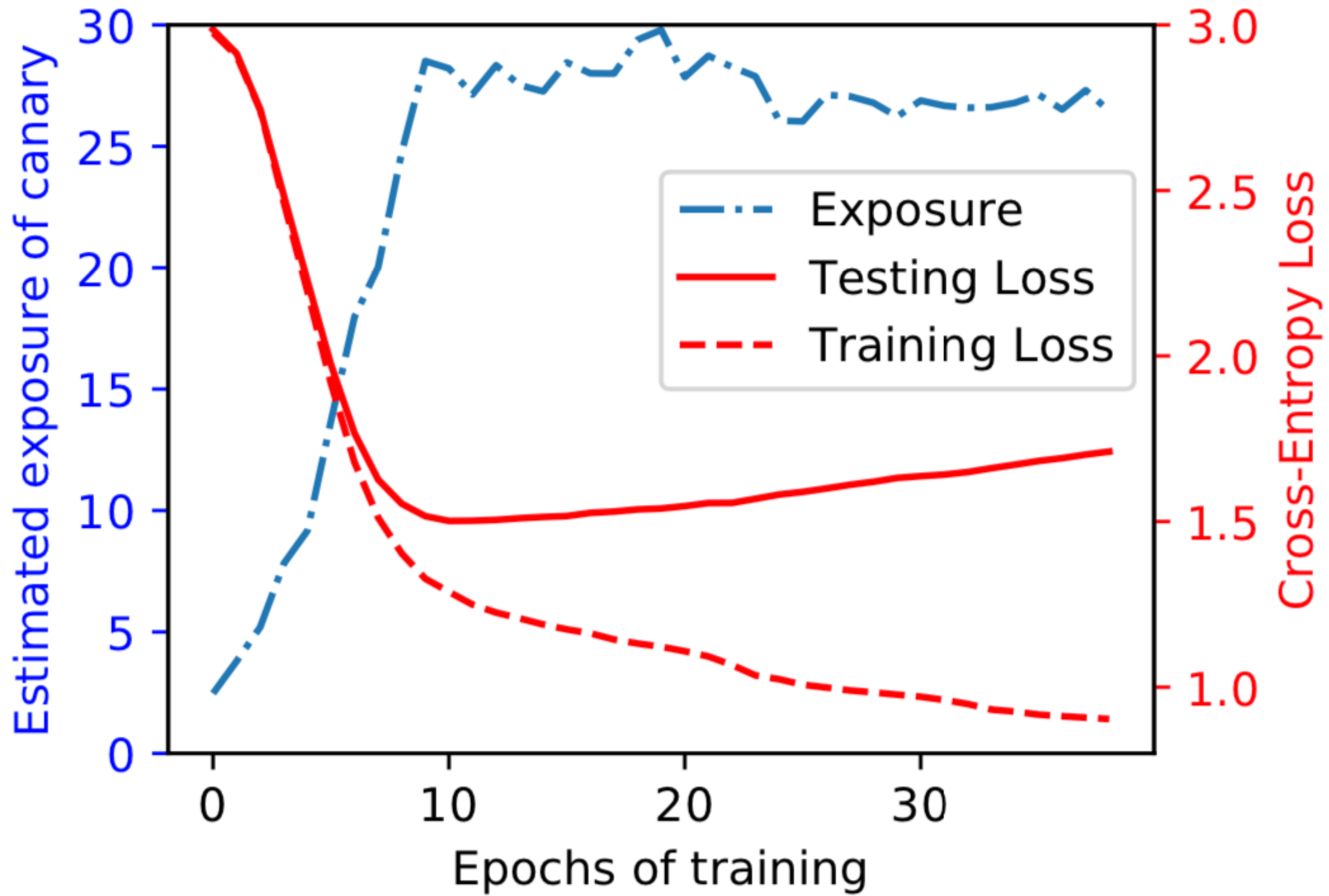
For the researcher:

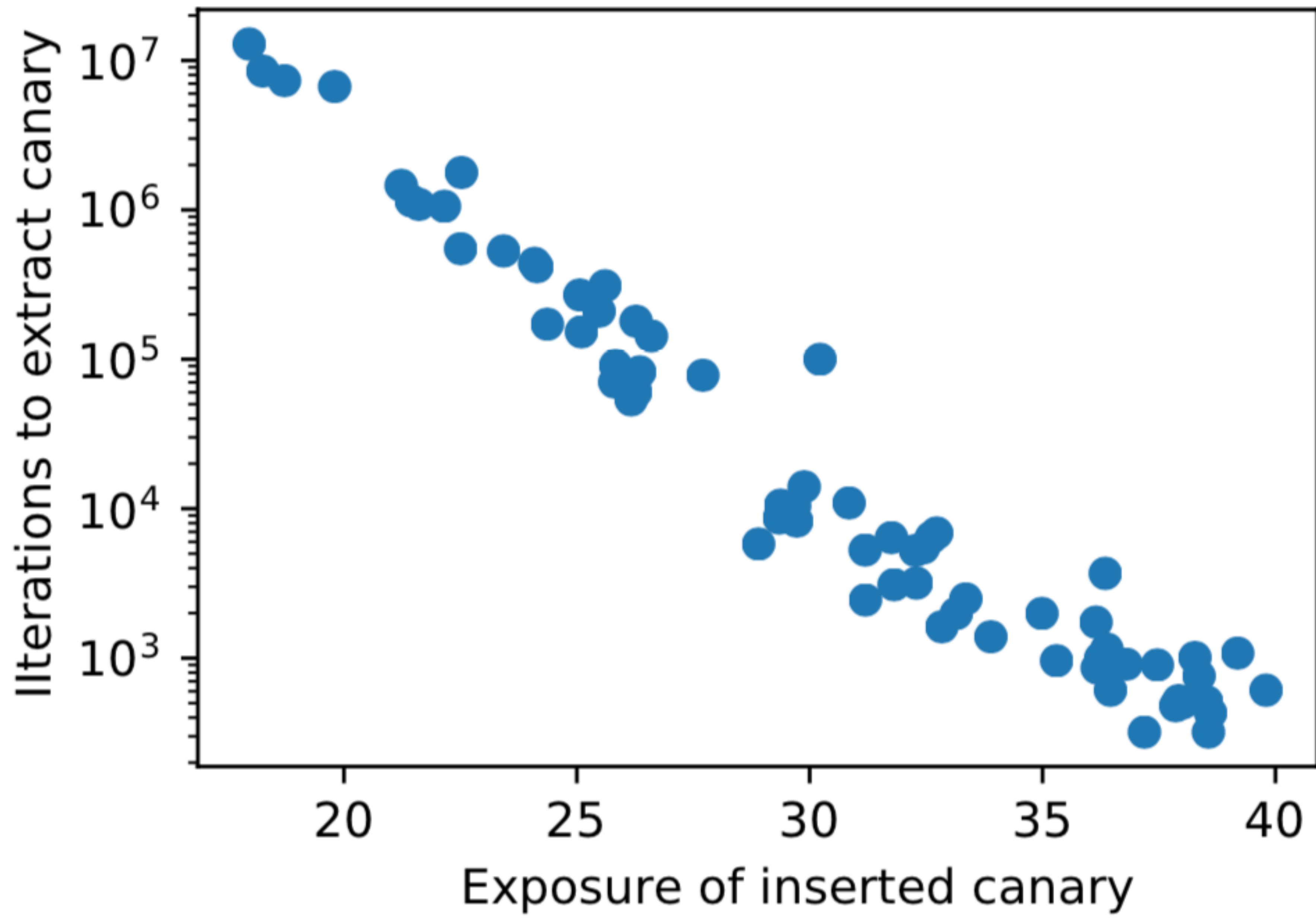
Measuring lower-bounds on
memorization is practical and useful.

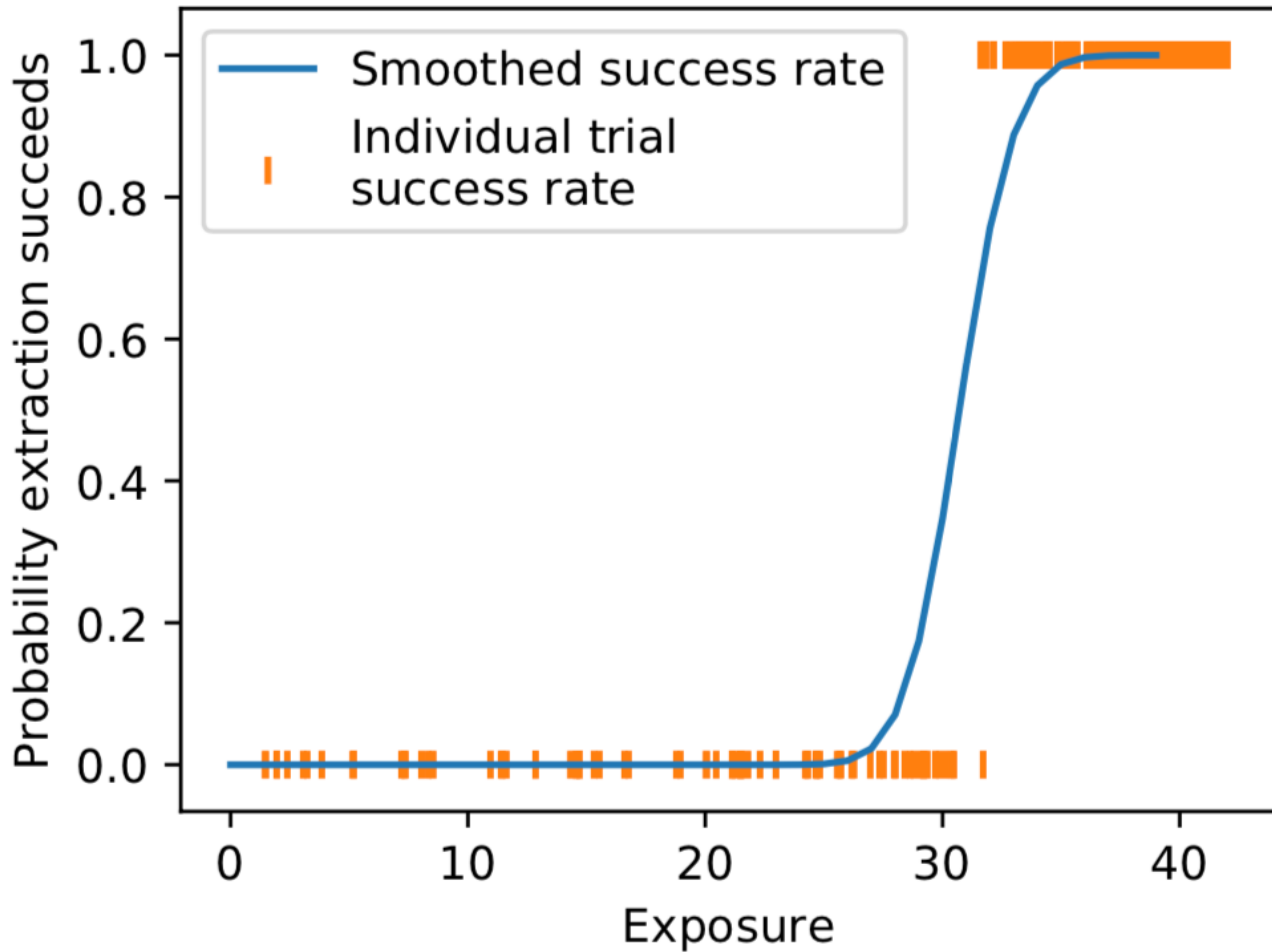
Questions

Backup Slides









User	Secret Type	Exposure	Extracted?
A	CCN	52	✓
B	SSN	13	
C	SSN	16	
	SSN	10	
	SSN	22	
D	SSN	32	✓
F	SSN	13	
G	CCN	36	
	CCN	29	
	CCN	48	✓

	Optimizer	ϵ	Test Loss	Estimated Exposure	Extraction Possible?
With DP	RMSProp	0.65	1.69	1.1	
	RMSProp	1.21	1.59	2.3	
	RMSProp	5.26	1.41	1.8	
	RMSProp	89	1.34	2.1	
	RMSProp	2×10^8	1.32	3.2	
	RMSProp	1×10^9	1.26	2.8	
	SGD	∞	2.11	3.6	
No DP	SGD	N/A	1.86	9.5	
	RMSProp	N/A	1.17	31.0	✓