# Security (and Privacy) in Machine Learning

Nicholas Carlini

*University of California, Berkeley*
*(now Google Brain)*

# This talk: neural networks

# The New York Times

## The Race for Self-Driving Cars

Autonomous cars have arrived. Major automakers have been investing billions in development, while tech players like Uber and Google's parent company have been testing their versions in American cities.

ALEX DAVIES TRANSPORTATION 03.13.18 12:15 PM

# WAYMO TAKES THE FINAL STEP BEFORE LAUNCHING ITS SELF-DRIVING CAR SERVICE

# Google AI defeats human Go champion
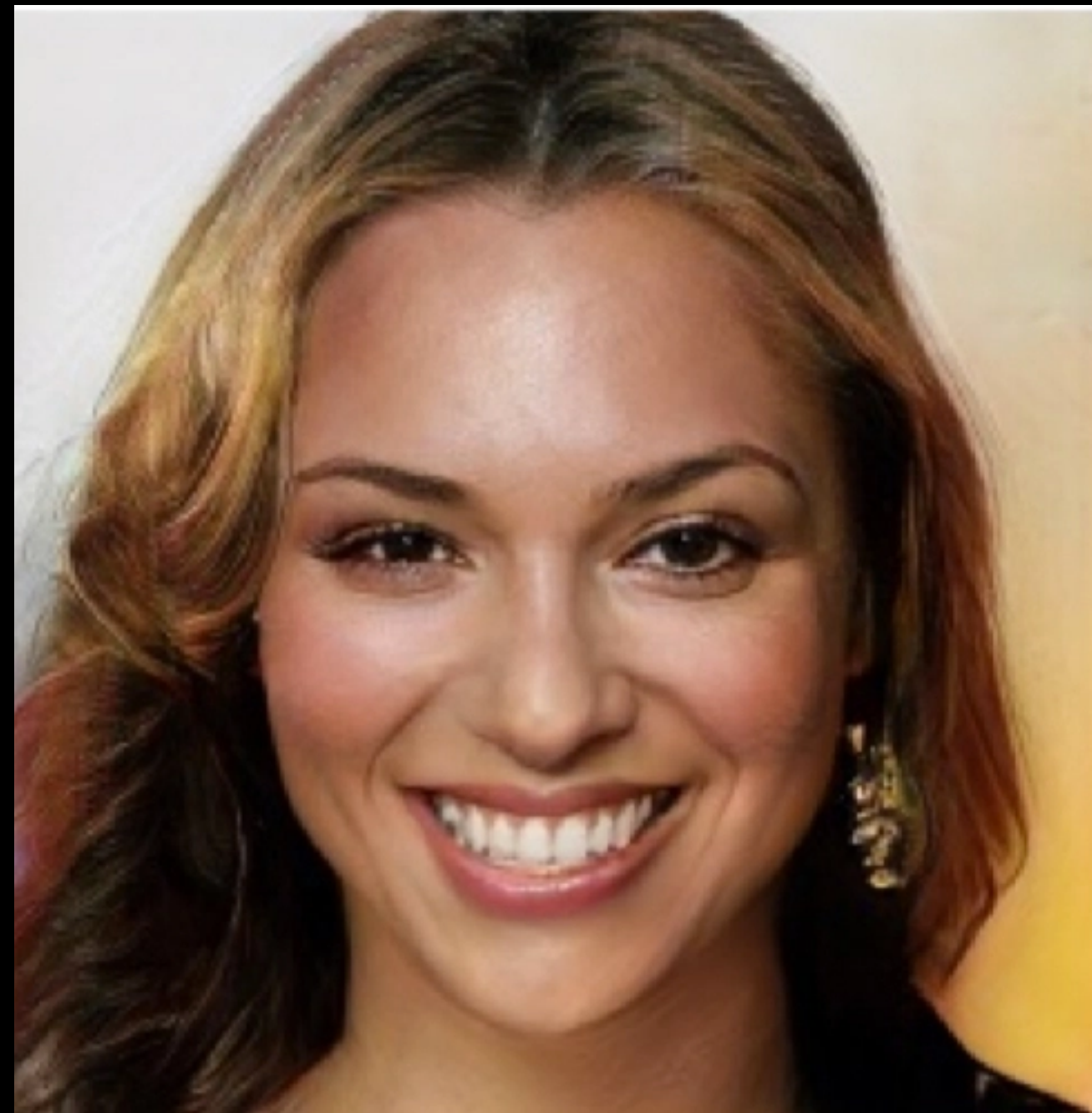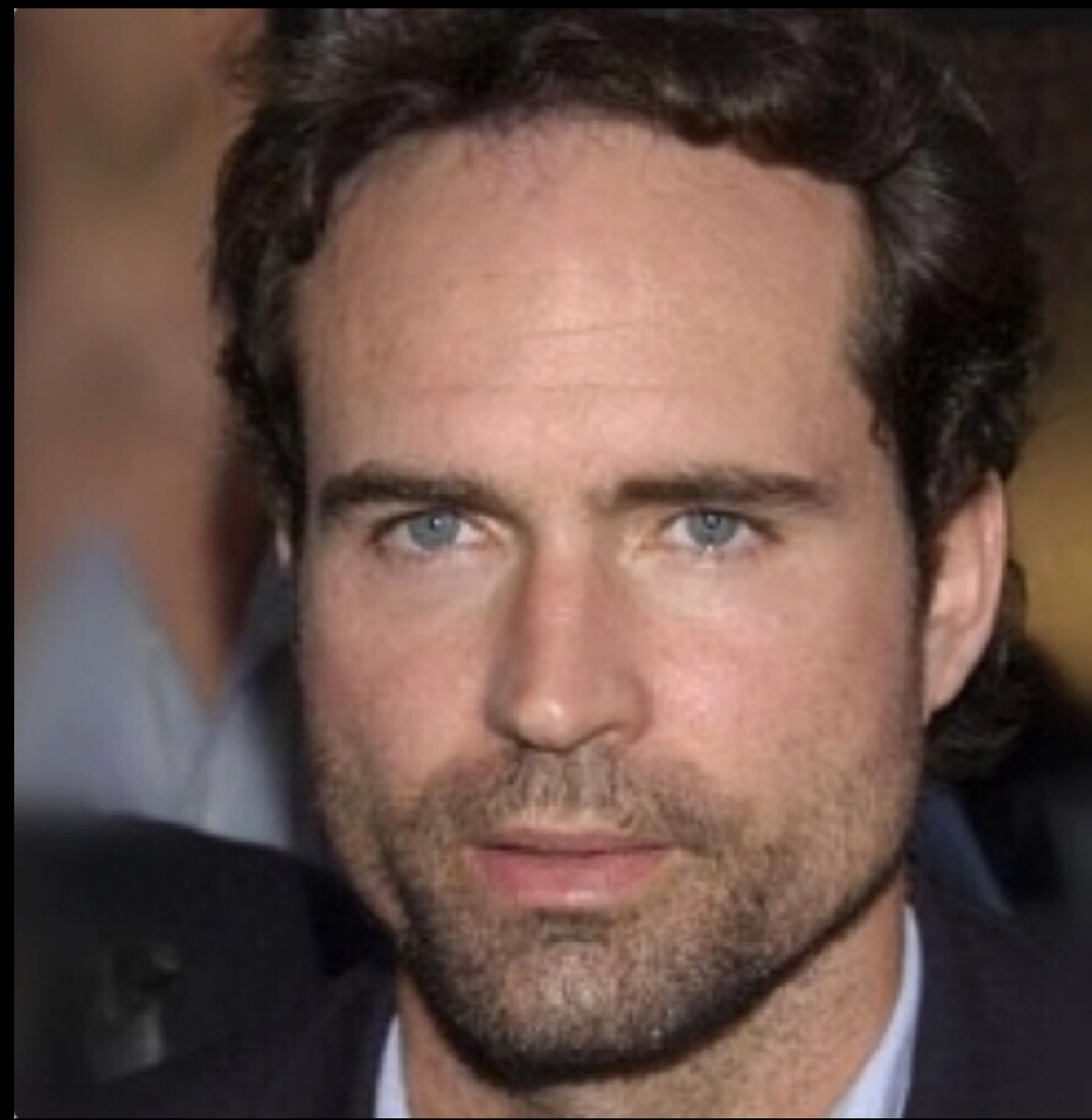
🕓 25 May 2017

≪ Share

---

**RESEARCH ARTICLE**

# Superhuman AI for heads-up no-limit poker: Libratus beats top professionals

**Noam Brown, Tuomas Sandholm**[*]

**+** See all authors and affiliations

*Science* 17 Dec 2017:
eaao1733
DOI: 10.1126/science.aao1733

What can I help you with?

amazon

# Machine learning is amazing

# But there's a catch

# Understandability

# This talk:

Discuss security & privacy problems being studied in the research community

# What this talk is *not*

$$U(x) = \left(\frac{1}{L}\sum_{i=1}^{L}\|F_r(x)\|\right) - \left\|\frac{1}{L}\sum_{i=1}^{L}F_r(x)\right\|$$

minimize $\mathcal{D}(x, x+\delta)$

such that $f(x+\delta) \leq 0$

$x + \delta \in [0,1]^n$

$$\Pr_{t\in\mathcal{R}}\left[\mathrm{Px}_\theta(s[t]) \leq \mathrm{Px}_\theta(s[r])\right]$$

$$= \sum_{v\leq \mathrm{Px}_\theta(s[r])}\Pr_{t\in\mathcal{R}}\left[\mathrm{Px}_\theta(s[t]) = v\right].$$

$$G(x)_i = \begin{cases} Z(x)_i & \text{if } i \leq N \\ (1+U(x)-\tau)\cdot \max_i Z(x)_i & \text{if } i = N+1 \end{cases}$$

$$\nabla_x f(g(x))|_{x=\hat{x}} \approx \nabla_x f(x)|_{x=g(\hat{x})}$$

$$\ell(\boldsymbol{x}) = \sum_i \max\left(\max_{t\in\{\epsilon, ""\}} f(\boldsymbol{x})_t^i - \max_{t'\notin\{\epsilon, ""\}} f(\boldsymbol{x})_{t'}^i, 0\right).$$

$$\Pr(\boldsymbol{p}|\boldsymbol{y}) = \sum_{\pi\in\Pi(\boldsymbol{p},\boldsymbol{y})}\Pr(\pi|\boldsymbol{y}) = \sum_{\pi\in\Pi(\boldsymbol{p},\boldsymbol{y})}\prod_i \boldsymbol{y}_{\pi^i}^i$$

# What this talk is *not*

**Definition 4.** *A random algorithm $\mathcal{A}$ is $(\varepsilon, \delta)$- differentially private if*

$$\mathbf{Pr}(\mathcal{A}(\mathcal{D})\in S) \leq \exp(\varepsilon)\cdot\mathbf{Pr}(\mathcal{A}(\mathcal{D}')\in S) + \delta$$

$$\alpha_{pq} = \sum_{i\in\{p,q\}}\frac{\partial Z(x)_t}{\partial x_i}$$

$$\beta_{pq} = \left(\sum_{i\in\{p,q\}}\sum_j \frac{\partial Z(x)_j}{\partial x_i}\right) - \alpha_{pq}$$

$$\theta^* = \arg\min_\theta \mathbb{E}_{x\in\mathcal{X}}\left[\max_{\delta\in[-\epsilon,\epsilon]^N}\ell(x+\delta; F_\theta)\right].$$

$$= -\log \mathbb{E}_{r'\in\mathcal{R}}\left[\mathbb{1}(L_\theta(s[r']) \leq L_\theta(s[r]))\right]$$

$$= -\log\left(1\cdot\frac{|\{r'\in\mathcal{R}: L_\theta(s[r'])\leq L_\theta(s[r])\}|}{|\mathcal{R}|}\right.$$

$$\left. 0\cdot\frac{|\mathcal{R}| - |\{r'\in\mathcal{R}: L_\theta(s[r'])\leq L_\theta(s[r])\}|}{|\mathcal{R}|}\right)$$

$$= -\log\left(\frac{|\{r'\in\mathcal{R}: L_\theta(s[r'])\leq L_\theta(s[r])\}|}{|\mathcal{R}|}\right)$$

$$= -\left(\log\mathbf{rank}_\theta(s[r]) - \log|\mathcal{R}|\right)$$

$$= \log|\mathcal{R}| - \log\mathbf{rank}_\theta(s[r])$$

What this talk *is*

What are the security problems
in machine learning today?

French Bulldog

(95%)

Old English Sheepdog

(83%)

Greater
Swiss
Mountain
Dog

(78%)

Siberian Husky

(81%)

Great Dane

(67%)

Beagle

(96%)

Guacamole

(99.99%)

Golden Retriever

(96%)

Guacamole

(99.99%)

# These phenomena are known as
# **adversarial examples**

B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Srndic, P. Laskov, G. Giacinto, and F. Roli. Evasion attacks against machine learning at test time. 2013.
C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. ICLR 2014.
I. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. 2014.

88% **tabby cat**   →   adversarial perturbation   →   99% **guacamole**

(a)  (b)  (c)

What does this have
to do with voice?

We use these same classification approaches for speech recognition.

# Attacks on Android, circa 2015

# State-of-the-art in 2015

It's been three years.

Can we do better?

# Feynman Algorithm

1. Write down the problem.
2. Think very hard.
3. Write down the answer.

# Towards Evaluating the Robustness of Neural Networks

Nicholas Carlini    David Wagner
University of California, Berkeley

# Audio Adversarial Examples: Targeted Attacks on Speech-to-Text

Nicholas Carlini    David Wagner
University of California, Berkeley

# Mozilla's DeepSpeech

# Mozilla's DeepSpeech transcribes this

Mozilla's DeepSpeech
transcribes this as

"most of them were staring
quietly at the big table"

[adversarial]

"It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity"

Why is this so much stealthier?

# It works on music, too

DeepSpeech transcribes
"speech can be embedded in music"

# And can "hide" speech

DeepSpeech does not hear any speech in this audio sample
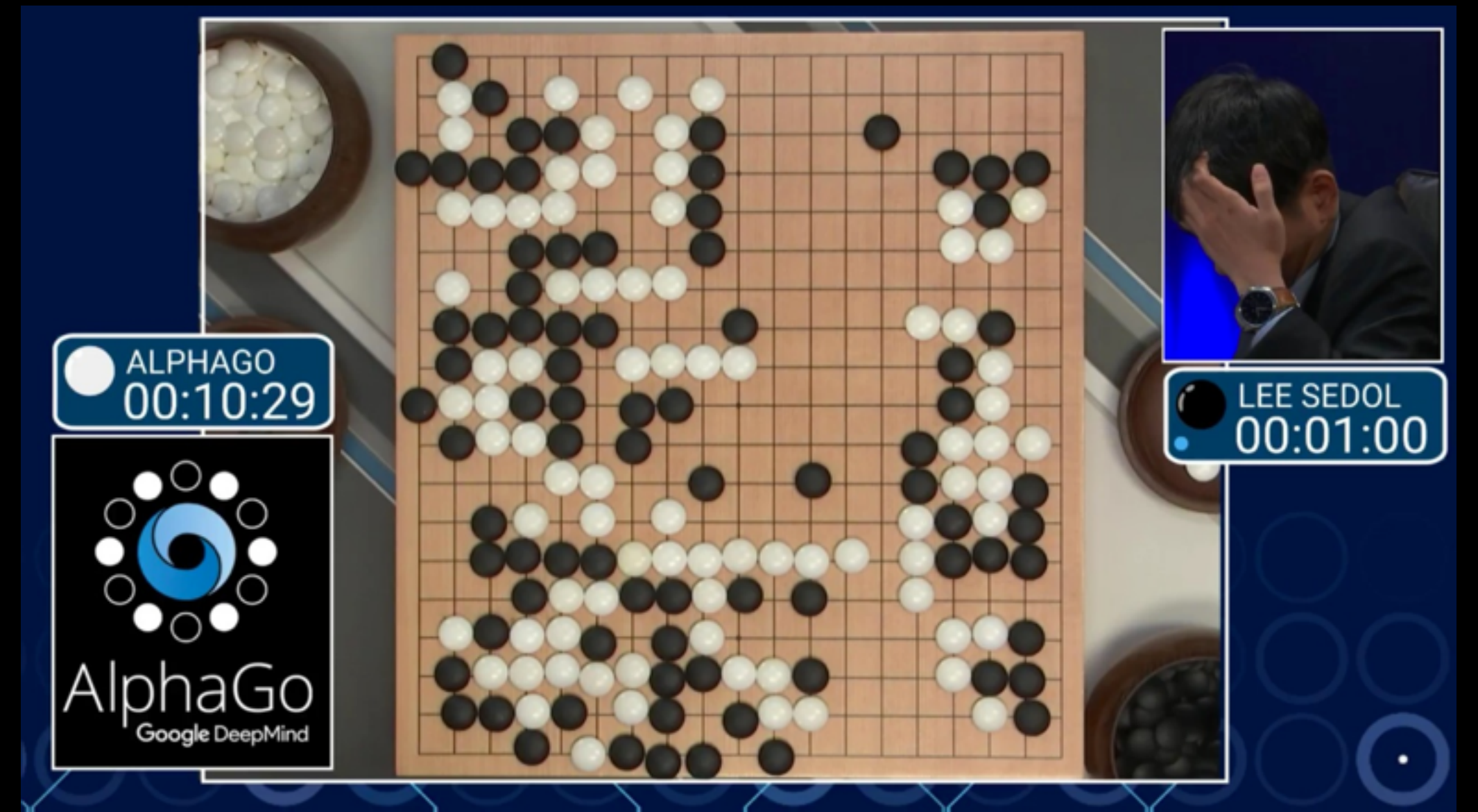
# That's a lot of problems

# Do you have any solutions?

# Sorry, no.

# This is an active area of research.

# Ask me again in two years.

Yes, machine
learning gives
**amazing** results

However, there are also significant **vulnerabilities**



Guacamole (99%)

# Questions?

More Details:
https://nicholas.carlini.com