# Current Research Thrusts

**Secure** AI    **Interpretable** AI

Why focus on them?
How are they related?

AI now used in safety-critical applications. Important to study threats & countermeasures.

Secure AI

The self-driving Uber was traveling north at about 40 m.p.h.

New York Times, 2018

How a Self-Driving Uber Killed a Pedestrian in Arizona

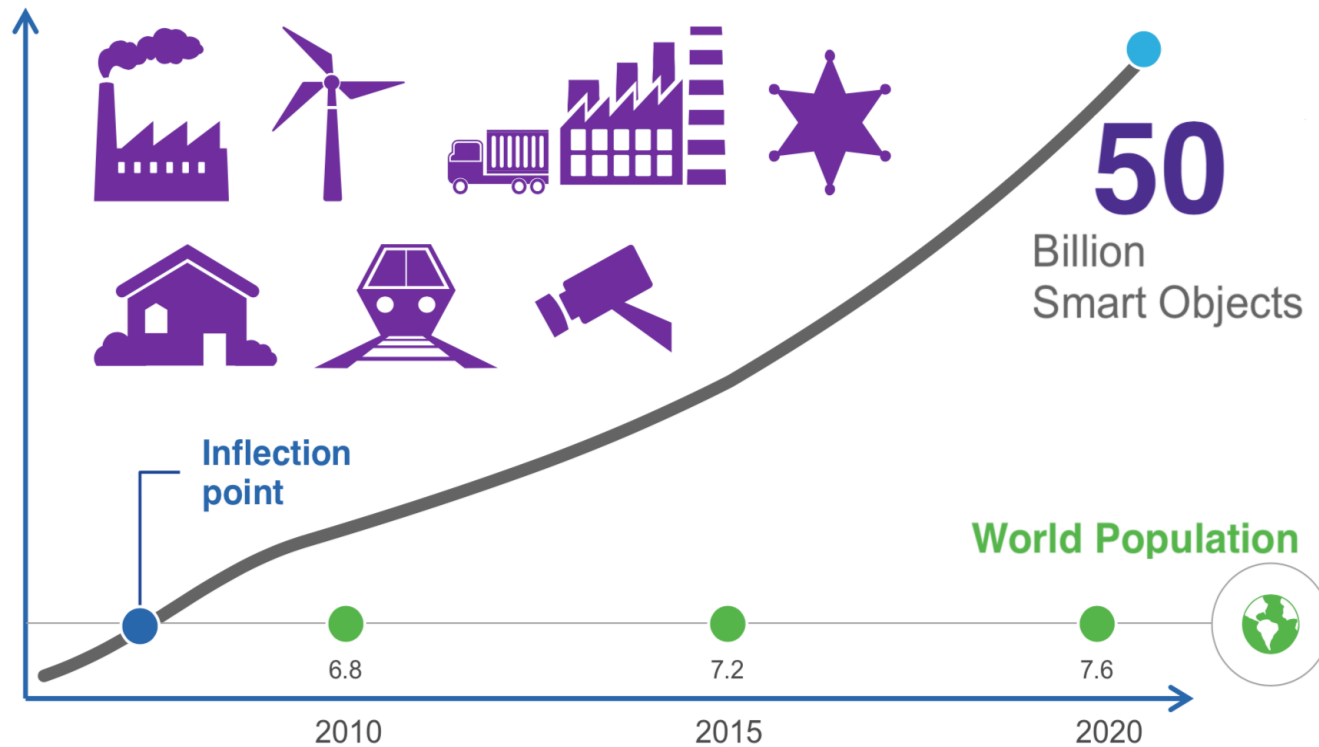# AI Security Problems Are Everywhere



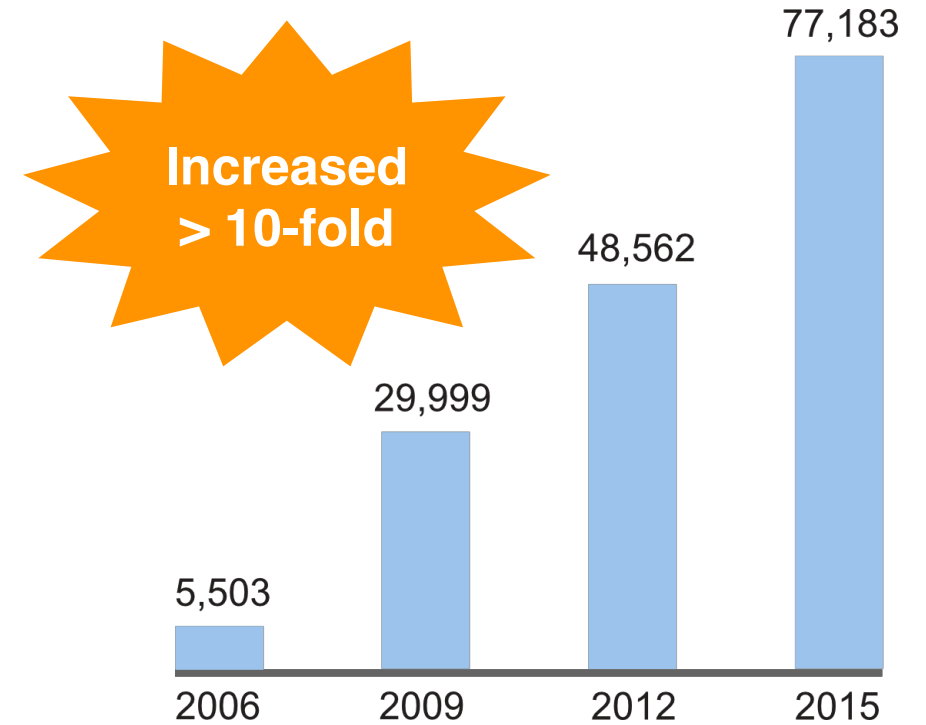"THE TOASTER HAS BEEN HACKED INTO THINKING IT'S A BLENDER."

klossner



Smart toaster does exist!
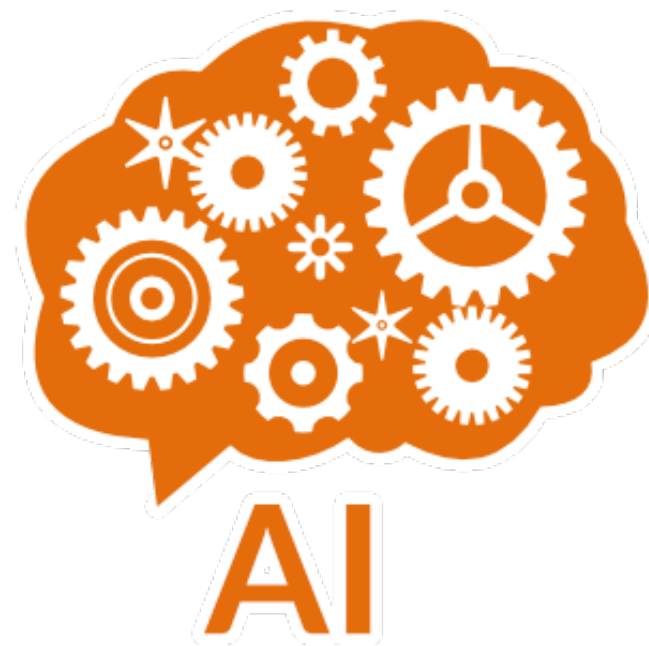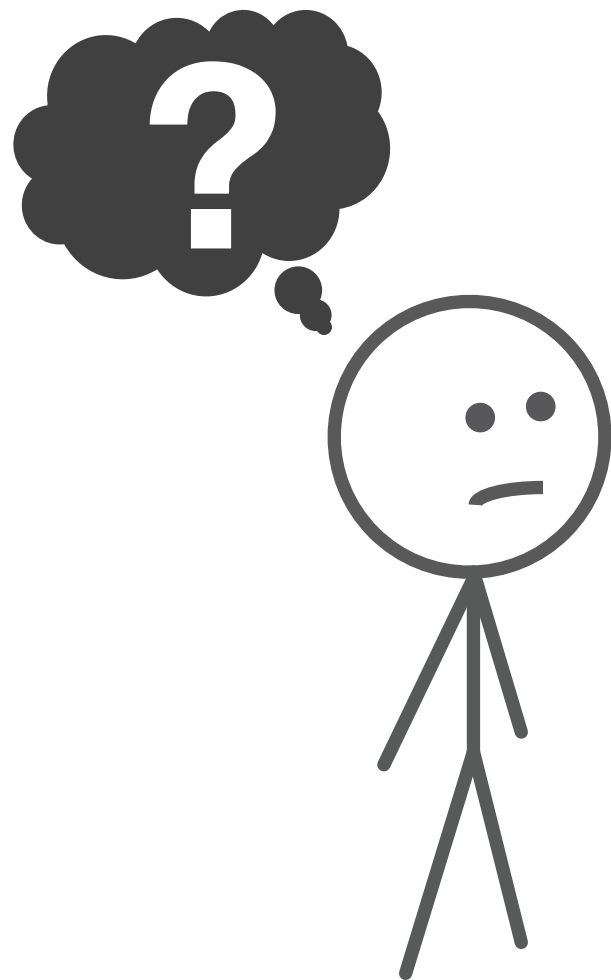
# AI Security is becoming increasingly important



50 Billion Smart Objects

Inflection point

World Population

6.8  2010
7.2  2015
7.6  2020

Source: Cisco

# incidents
reported by U.S. federal agencies

Increased > 10-fold

5,503 — 2006
29,999 — 2009
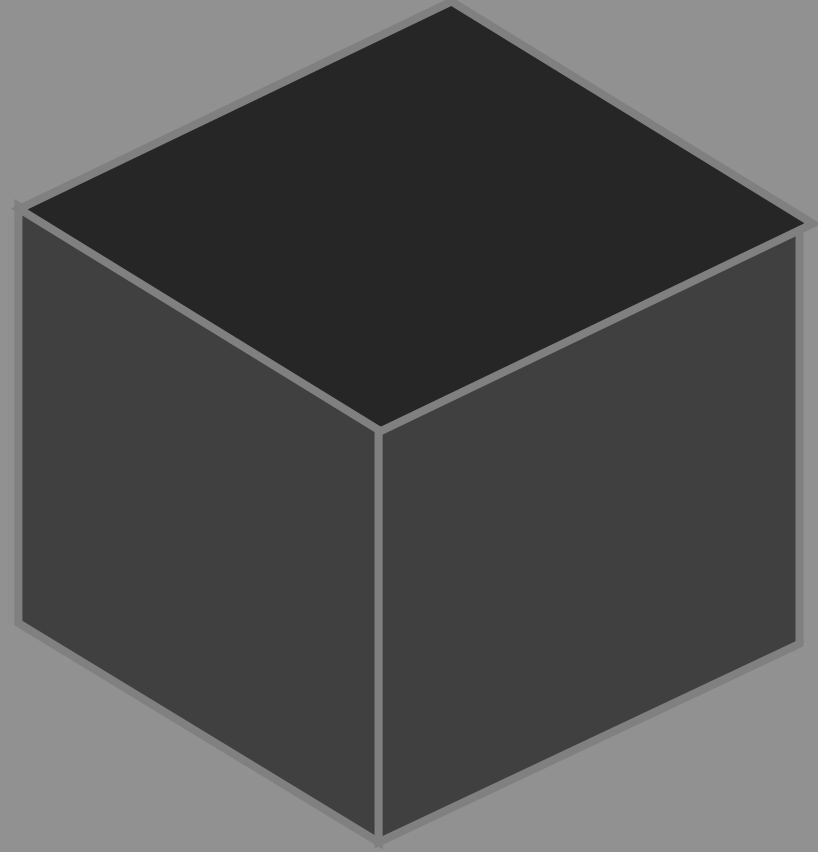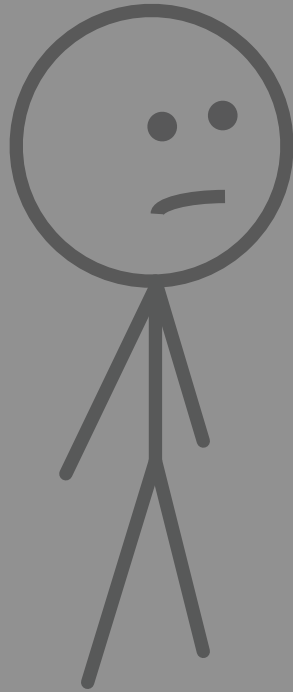48,562 — 2012
77,183 — 2015
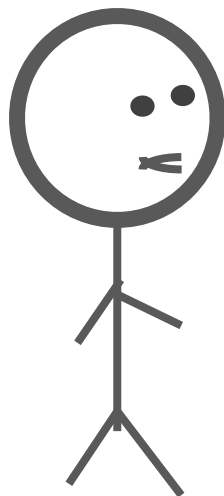
Source: US Department of Homeland Security

# How do we know if a defense for AI is working?

# AI models often used as black-box

# Interpretable AI

# Interpretable AI

Via **scalable, interactive, usable interfaces** to help people understand complex, large-scale ML systems.

**Secure AI**

**ShapeShifter** First attack fooling object detectors

**SHIELD & Adagio** Real-time defense

**Interpretable AI**

**Summit** Scalable interpretation for deep learning

**GAN Lab & CNN Explainer** Interactive learning

**Surveys** AI guidelines, visual analytics for deep learning
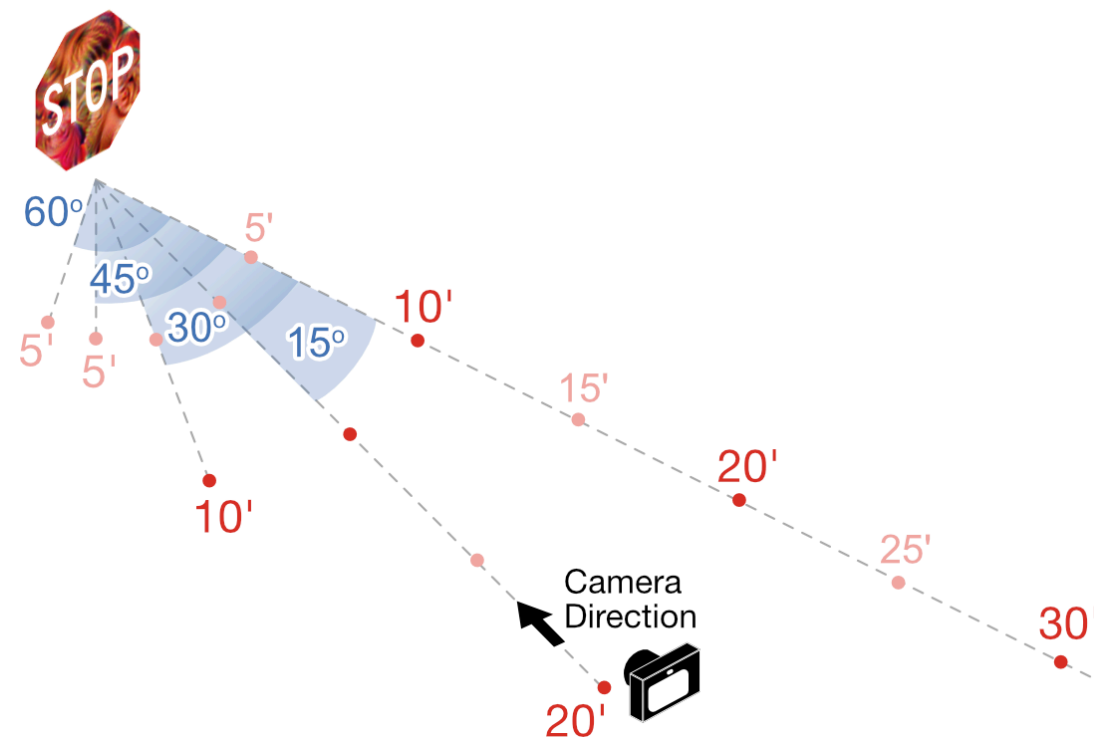
# Challenges of Physically Attacking Faster R-CNN
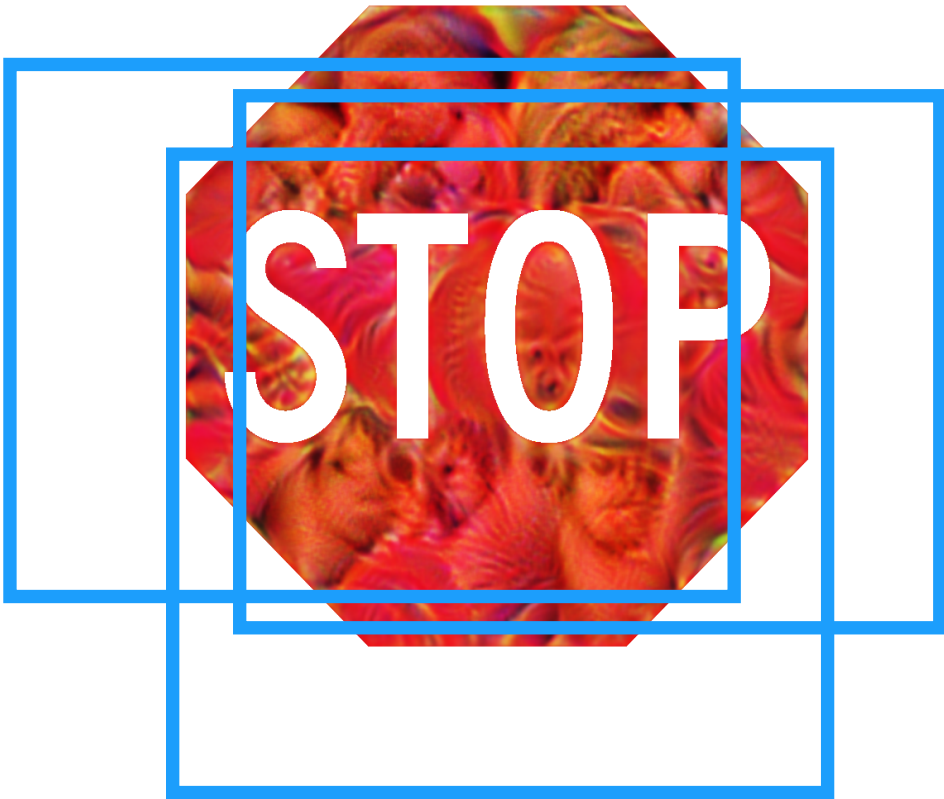
## 1. Multiple region proposals



## 2. Distances, angles, lightings

# Our Solution: **Fool Multiple Region Proposals**

Minimize: sum of classification losses + deviation loss



STOP ≈ STOP

Only perturb **RED** area
Human eye is less sensitive
to changes in darker color

# Our Solution: **Robust to Real-World Distortions**

Adapt Expectation over Transformation [Athalye et al, ICML'18]



Optimize over different backgrounds, scales, rotations, lightings

# ShapeShifter Motivates
# DARPA Program GARD (Defense for AI)
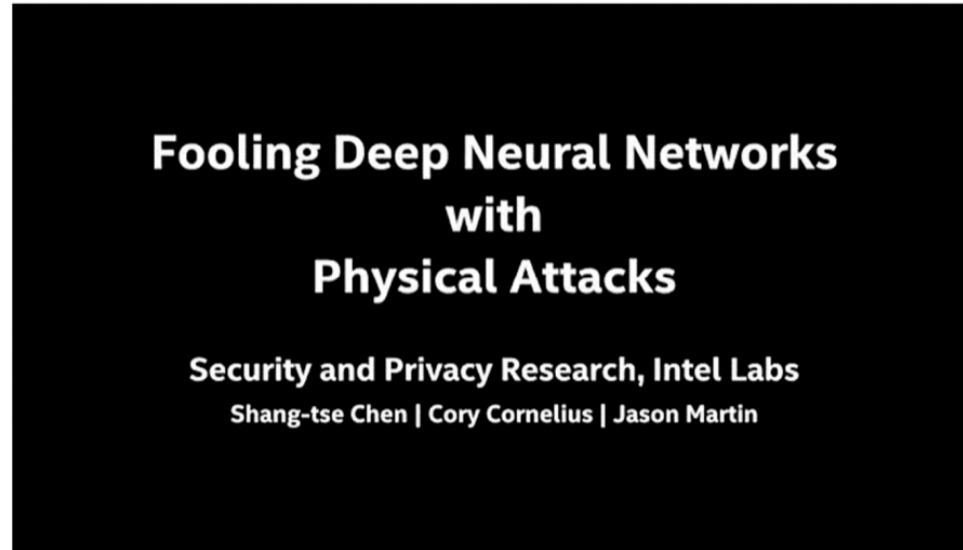


**State of the art: few physical attacks**

Graffiti:
(Evtimov et al., UC Berkeley, 2017)

Patch:
(Brown et al., Google, 2017)

3D Printed Objects:
(Athalye et al., MIT, 2017)

**Fooling Deep Neural Networks with Physical Attacks**

Security and Privacy Research, Intel Labs

Shang-tse Chen | Cory Cornelius | Jason Martin

(Intel / GTECH 2018)

• All physical attacks to date are White Box
• No current consideration of resource constraints

Highlights **ShapeShifter** as the state-of-the-art physical attack

https://www.darpa.mil/attachments/GARD_ProposersDay.pdf

# Adversarial Machine Learning Landscape

Attack

Defense

Our Focus:
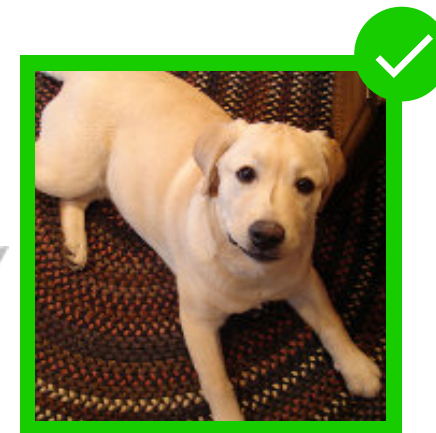**Fast & Practical**
(digital)

SHIELD
Secure Heterogeneous Image Ensemble with Localized Denoising

"Chain Mail" (Attacked)

Labrador Retriever

Real-time Compression Preprocessing

Vaccinated Deep Neural Network Ensemble

Correctly Classified

Correctly Classified

# SHIELD leverages JPEG compression



JPEQ Quality 80

JPEQ Quality 60

JPEQ Quality 40

JPEQ Quality 20

SHIELD's SLQ applies JPEG compression
of a random quality to each
8 x 8 block of the image

* larger blocks shown for presentation

# Defense Runtime Comparison
## (in seconds; shorter is better)

| Total Variation Denoising — Weight | | |
|---|---|---|
| 10 | 2049 | |
| 20 | 2041 | |
| 30 | 1743 | |
| 40 | 1723 | **>22x** Slower than JPEG-20 |

| Median Filter — Window Size | | |
|---|---|---|
| 5 | 3178 | |
| 3 | 1102 | **>14x** Slower than JPEG-20 |

| JPEG — Quality | | |
|---|---|---|
| 80 | 107 | |
| 60 | 92 | |
| 40 | 85 | |
| 20 | 77 | |

tested on 50,000 images from the ImageNet validation set

# ADAGIO

Interactive Experimentation with Adversarial Attack & Defense for Audio

- 🗣️ Upload your own audio sample
- ⚔️ Perform audio adversarial attack
- 📒 Apply compression to defend
- ▶️ Play audio, listen for differences



**ADAGIO** = **A**ttack & **D**efense for **A**udio in a **G**adget with **I**nteractive **O**perations

SUMMIT

MODEL
InceptionV1

DATASET
ImageNet

CLASSES
1,000

INSTANCES
1,281,024

LAYER
mixed

| 3a | 3b | 4a | 4b | 4c | 4d | 4e | 5a | 5b |

CLASS
white_wolf

INSTANCES
1299

ACCURACY
81.8%

PROBABILITIES

FILTER GRAPH

ADJUST WIDTH

ADJUST HEIGHT

• timber wolf

•white wolf
• malamute

• pembroke

• samoyed

• arctic fox
• shetland sheepdog
• lesser panda
• papillon    • keeshond
• collie
• chow

🔍 tench

**tench**                                   1.8%

🗋 red wolf                   69.9%

🗋 timber wolf               64.2%

🗋 arctic fox                87.1%

🗋 lion                      87.1%

🗋 chow                      87.1%

🗋 rottweiler               76.6%

🗋 silky terrier            63.3%

# Generative Adversarial Networks (GANs)

*"the most interesting idea in the last 10 years in ML"*

\- Yann LeCun



Face images generated by BEGAN [Berthelot et al., 2017]

# Why GANs are hard?

A GAN uses two *competing* neural networks

**Generator**
synthesizes outputs

**Discriminator**
spots fake

**Counterfeiter**
makes fake bills

**Police**
spots fake bills

# **CNN Explainer** also went viral! Try at **bit.ly/cnn-explainer**

⭐ 5.3K GitHub Stars    ❤️ 700 Likes    36K visitors, 151 countries