# Garbage In, Garbage Out

## *How purportedly great ML models can be screwed up by bad data*

**Hillary Sanders**
Data Scientist - operations team lead

SOPHOS

# What I'll show...

1. Model accuracy claimed by security ML researchers is misleading

2. It's generally biased in an overly optimistic direction

3. → Estimating the severity of that bias is important, and will help you make sure that your model isn't… garbage.

# Machine Learning

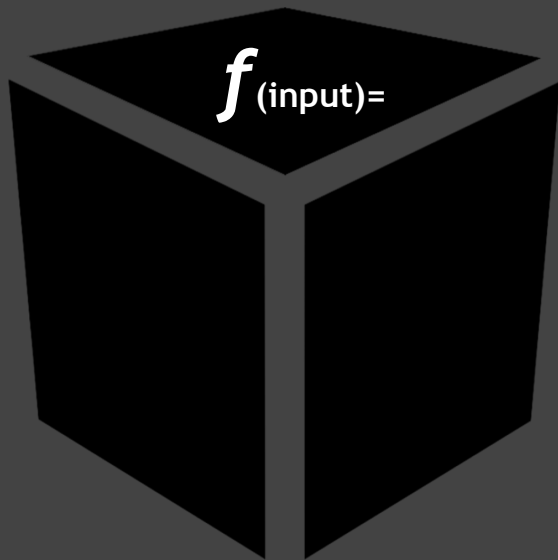$$f(\text{input}) = \text{output}$$

$$f(\text{http://www.trustus.evil.ru/paypal/login/}) = .944780$$

$$f(\text{https://www.facebook.com/}) = .019367$$

# Machine Learning

input → $f_{(input)}=$ → output

http://www.trustus.evil.ru/paypal/login/    .944780
http://gsbyntwqmem.mrjz5viern.ru/start_page.exe    .99981683
https://www.facebook.com/    .019367
http://imgur.com/r/cats/omgn4Zv    .008448

# *TRAINING*

# (Supervised) Machine Learning

Training Data

http://www.avit.ru/raypoa/bqjH/, 1
http://galovh/ru/start_page.exe, 1
https://www.facebook.com/, 0
http://imgur.com/r/cats/omgn4Zv, 0

Test Data

bwpqbyfykfizylpszfize.biz/, 1
http://git.demo.nick.net.nz/, 0

# (Supervised) Machine Learning

Training Data

http://www.evil.ru/payload/login/, 1
http://galoyh.ru/start_page.exe, 1
https://www.facebook.com/, 0
http://imgur.com/r/cats/omgn4Zv, 0

Test Data

0wprgbyfykfzkjlbydfzz.biz/, 1
http://git.demo.nick.net.nz/, 0

"**fitted**" model $f$

$f$ (input)

Training accuracy
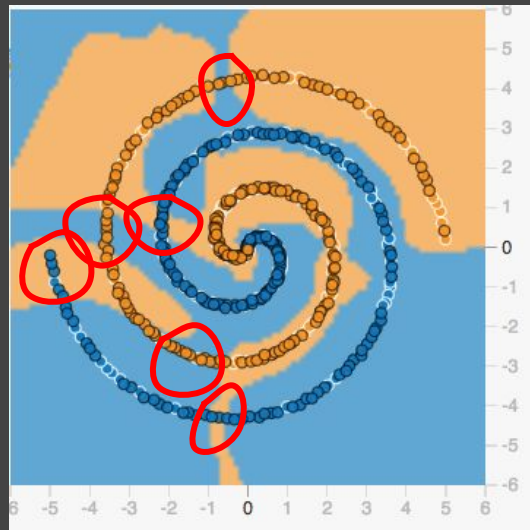
Test accuracy

# Training accuracy

# Training accuracy

# *Test* accuracy

# (Supervised) Machine Learning

Training Data

http://www.evil.ru/paypal/login/, 1
http://gdzyn.ru/start_page.exe, 1
https://www.facebook.com/, 0
http://imgur.com/r/cats/omgn4Zv, 0

Test Data

bwpqbyfykrlzylpudlze.biz/, 1
http://git.demo.nick.net.nz/, 0

"**fitted**" model *f*

*f* (input)

# (Supervised) Machine Learning

Training Data

http://www.evil.ru/payload/login/, **1**
http://gabryhn.ru/start_page.exe, **1**
https://www.facebook.com/, **0**
http://imgur.com/r/cats/omgn4Zv, **0**

Test Data

bwprjbyfyklrlzy/pudlzs.biz/, 1
http://git.demo.nick.net.nz/, 0

"**fitted**" model *f*

*f* (input)

**Test Accuracy**
.99 (error ≈ .01)
.01 (error ≈ .01)

# DEPLOYMENT
## accuracy

# (Supervised) Machine Learning

**Training Data**

http://www.evil.ru/payload/loghi/, **1**
http://galoyni.ru/start_page.exe, **1**
https://www.facebook.com/, **0**
http://imgur.com/r/cats/omgn4Zv, **0**

"**fitted**" model *f*

*f* (input)

**Test Data**

bwprgbyfyfkrlzyljxu6zz.biz/, 1
http://git.demo.nick.net.nz/, 0

**Deployment Data!**

???.evil.com
???.good.com

?

**Deployment Accuracy?**
**???, ???**

|            | Training | Testing |
|------------|----------|---------|
| **Lab**    | ✅       | ✅      |
| *Deployment* | 🟠 ? | 🟠 ? |

|  | Training | Testing |
|---|---|---|
| Lab | ✅ | ✅ |
| *Deployment* | ❌ | ❌ |

|  | | Training | Testing |
|---|---|---|---|
| | Lab | ✅ | ✅ |
| | *Deployment* | ❌ | ❌ |

|  | Training | Testing |
|---|---|---|
| Lab | ✅ | ✅ |
| *Deployment* | ❌ | ❌ |

|  | | Training | Testing |
|---|---|---|---|
| Lab | | ✓ | ✓ |
| *Deployment* | | ✗ | ✗ |

# Train / Test "Sensitivity Analysis": *Identifying training data that leads to improved and consistent performance on new datasets*

## Train and test the same model across different datasets, and evaluate the results:

1. What training datasets generalize better to others?

2. How sensitive is a model's accuracy to changes in test datasets?

# Train / Test "Sensitivity Analysis"
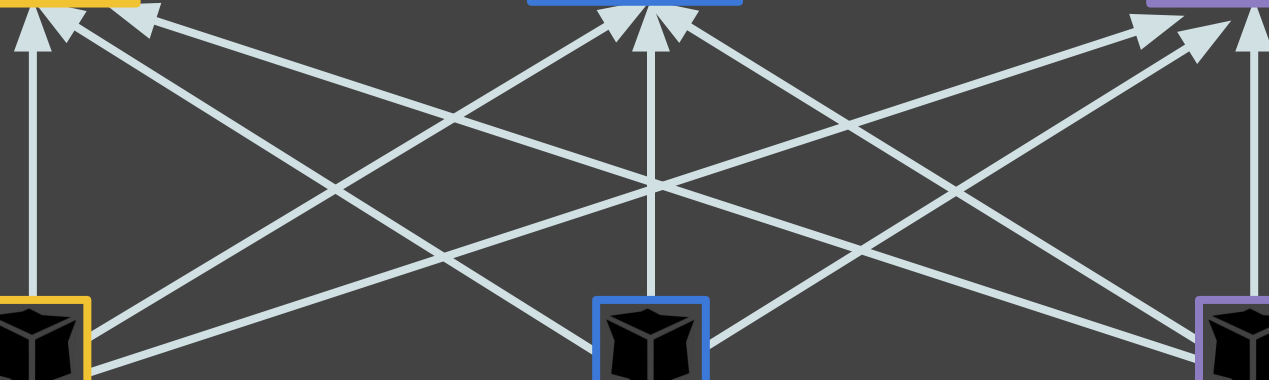
# IRL!

# Train / Test "Sensitivity Analysis"

1. Model Used

2. Accuracy Metric Used: AUC
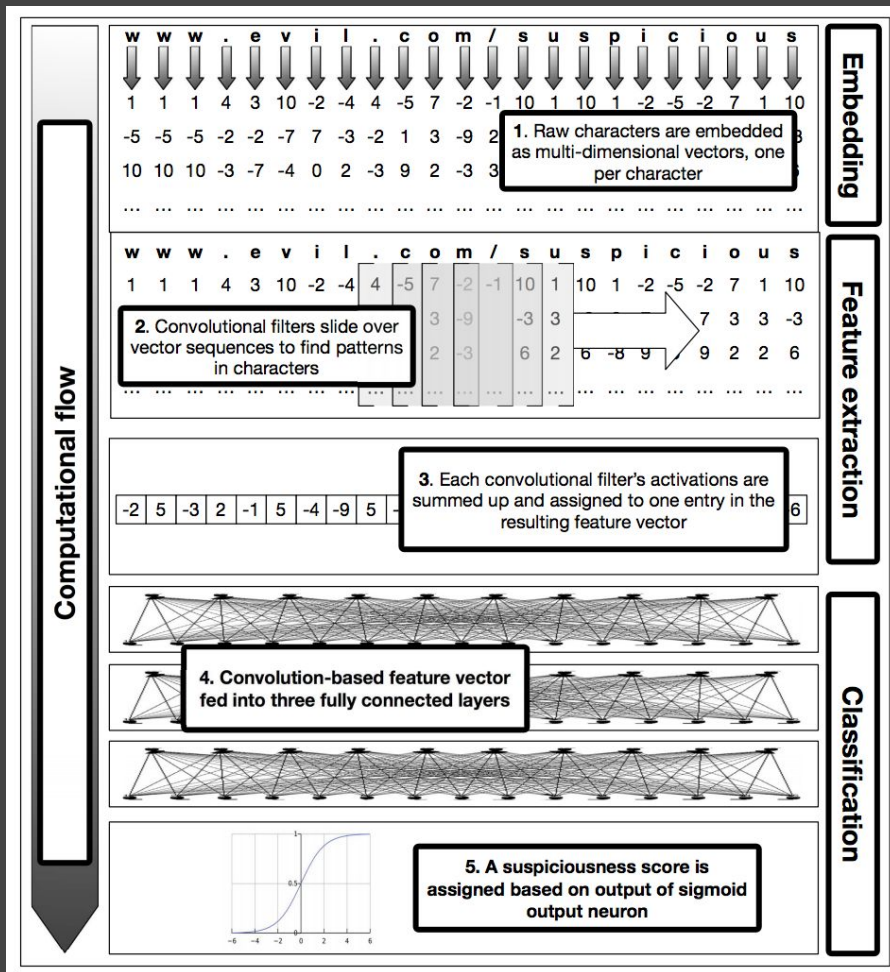
3. Datasets Used

4. Results!

# Train / Test "Sensitivity Analysis"

1. **Model Used**

2. Accuracy Metric Used: AUC

3. Datasets Used

4. Results!

# URL Model

*A Character-Level Convolutional Neural Network with Embeddings For Detecting Malicious URLs, File Paths and Registry Keys*
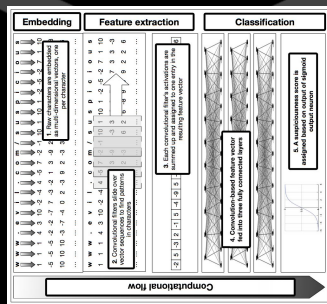
*Joshua Saxe, Konstantin Berlin*
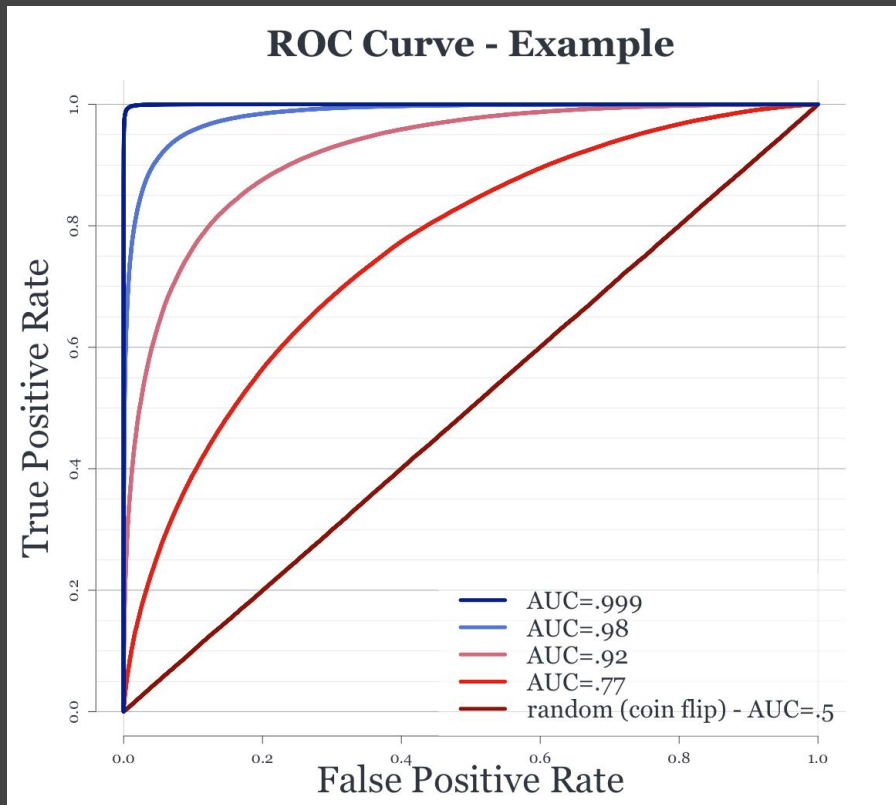
# Train / Test "Sensitivity Analysis"

1. Model Used

2. Accuracy Metric Used: AUC

3. Datasets Used

4. Results!

# AUC = "Area Under the [ROC] Curve"

# Train / Test "Sensitivity Analysis"

1. Model Used

2. Accuracy Metric Used: AUC

3. Datasets Used

4. Results!

## CommonCrawl & PhishTank

10 million URLs from January 2017*

≈ **20k** malware samples

*\* plus pre-Jan '17 phishtank malicious URLs, due to lack of data*

## Sophos

10 million internal URLs from January 2017
≈ 4% malware

## VirusTotal

10 million URLs from January 2017

≈ 4% malware

# CommonCrawl & Phishtank

## Sophos

## VirusTotal

**10 million URLs from January 2017***

**≈ 20k malware samples**

*\* plus pre-Jan '17 phishtank malicious URLs, due to lack of data*

**10 million internal URLs from January 2017**
**≈ 4% malware**

**10 million URLs from January 2017**

**≈ 4% malware**

## CommonCrawl & Phishtank

10 million URLs from January 2017*

≈ 20k malware samples

*plus pre-Jan '17 phishtank malicious URLs, due to lack of data*

## Sophos

10 million internal URLs from January 2017
≈ 4% malware

## VirusTotal

10 million URLs from January 2017

≈ 4% malware

## CommonCrawl & Phishtank

10 million URLs from January 2017*

≈ 20k malware samples

*plus pre-Jan '17 phishtank malicious URLs, due to lack of data*

## Sophos

10 million internal URLs from January 2017

≈ 4% malware

## VirusTotal

10 million URLs from January 2017

≈ 4% malware

**CommonCrawl & PT model**

*(January '17 data)*

**Sophos model**

*(January '17 data)*
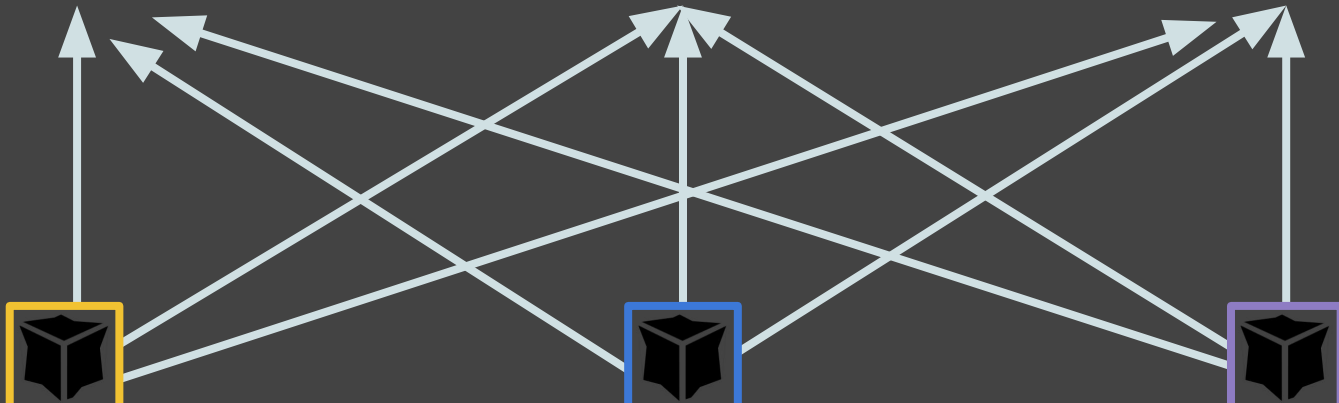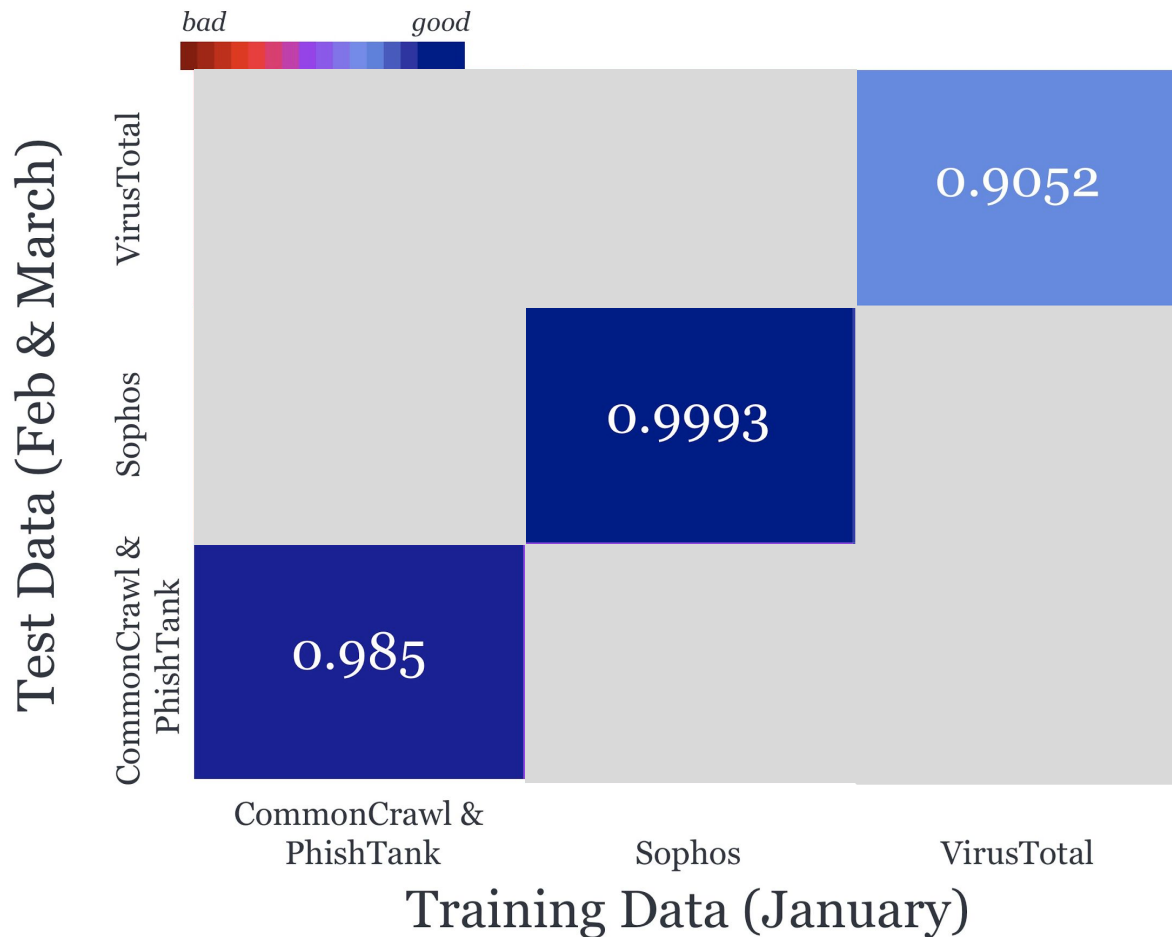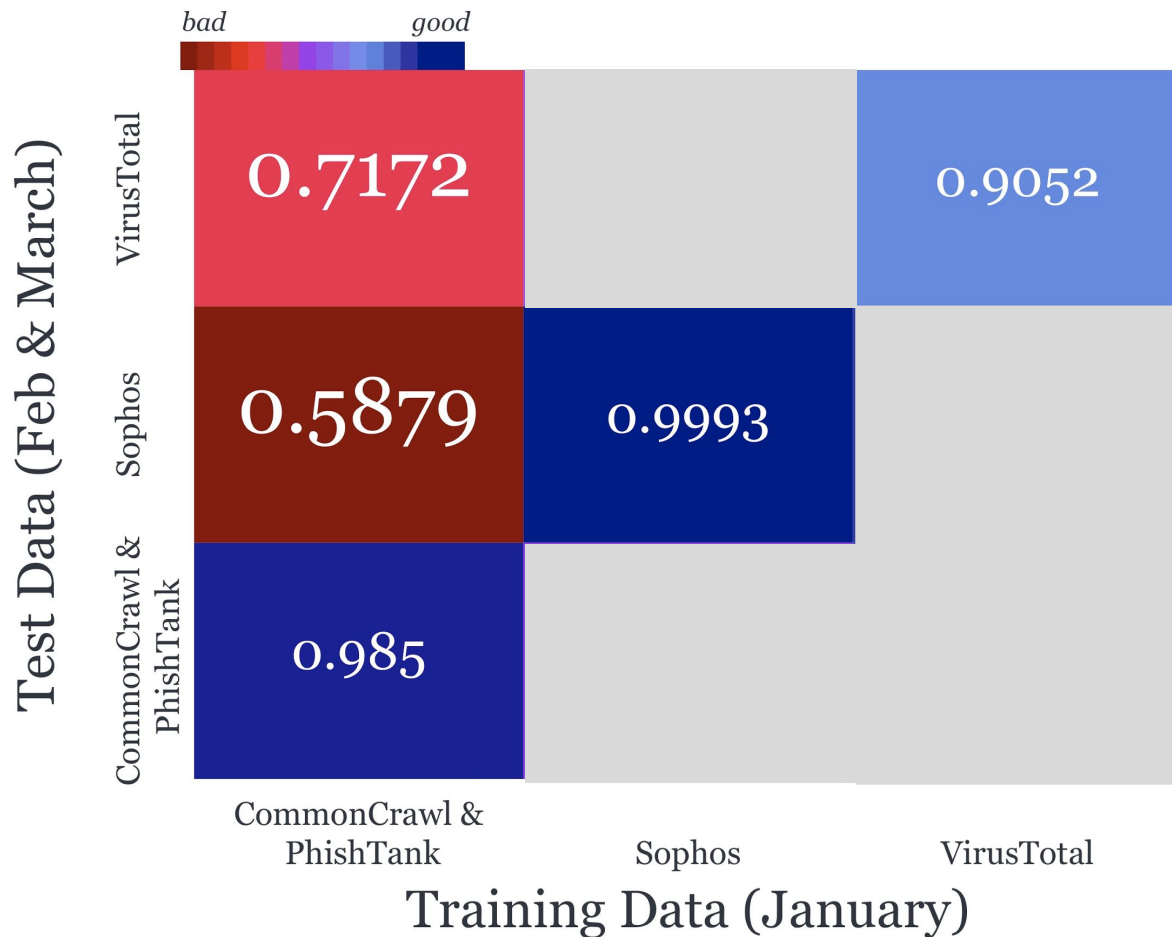
**VirusTotal model**

*(January '17 data)*

# Train / Test "Sensitivity Analysis"

1. Model Used

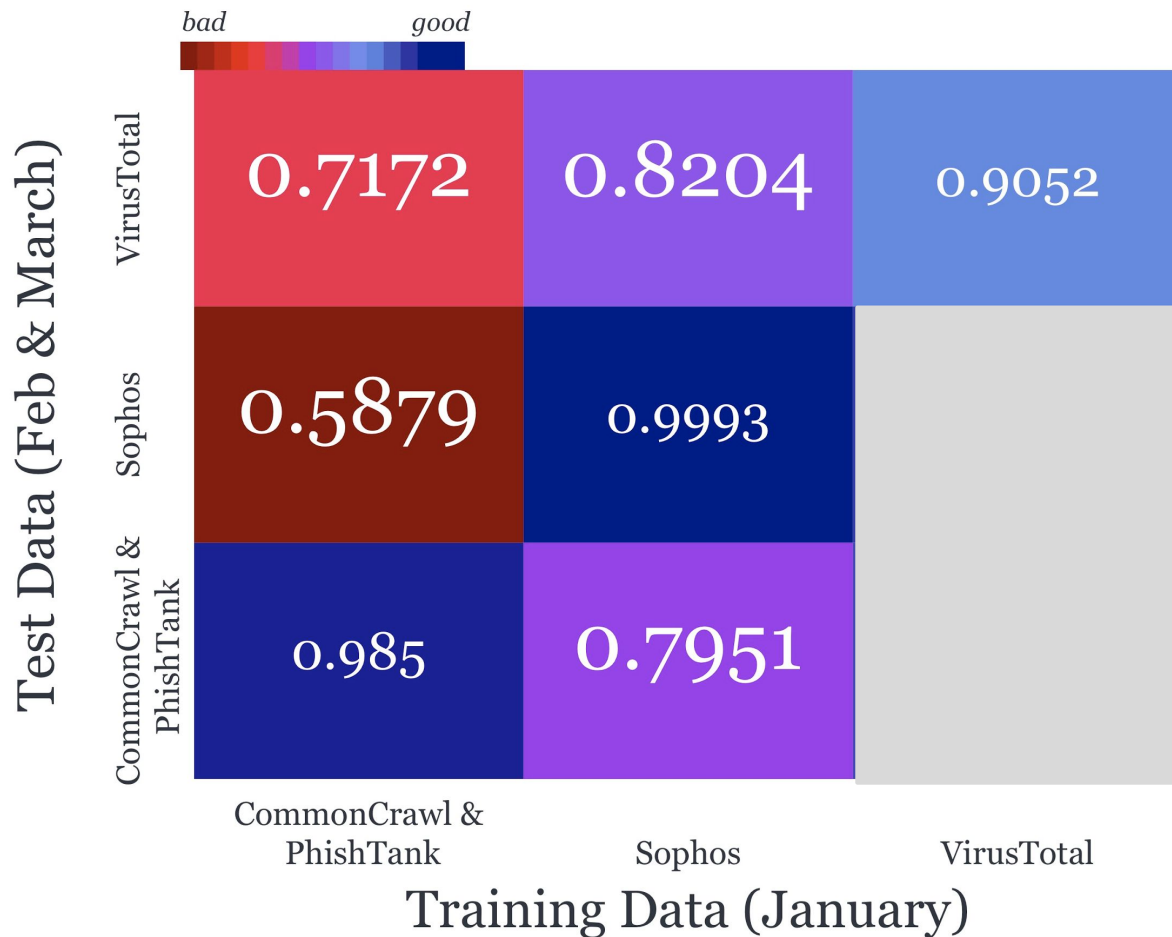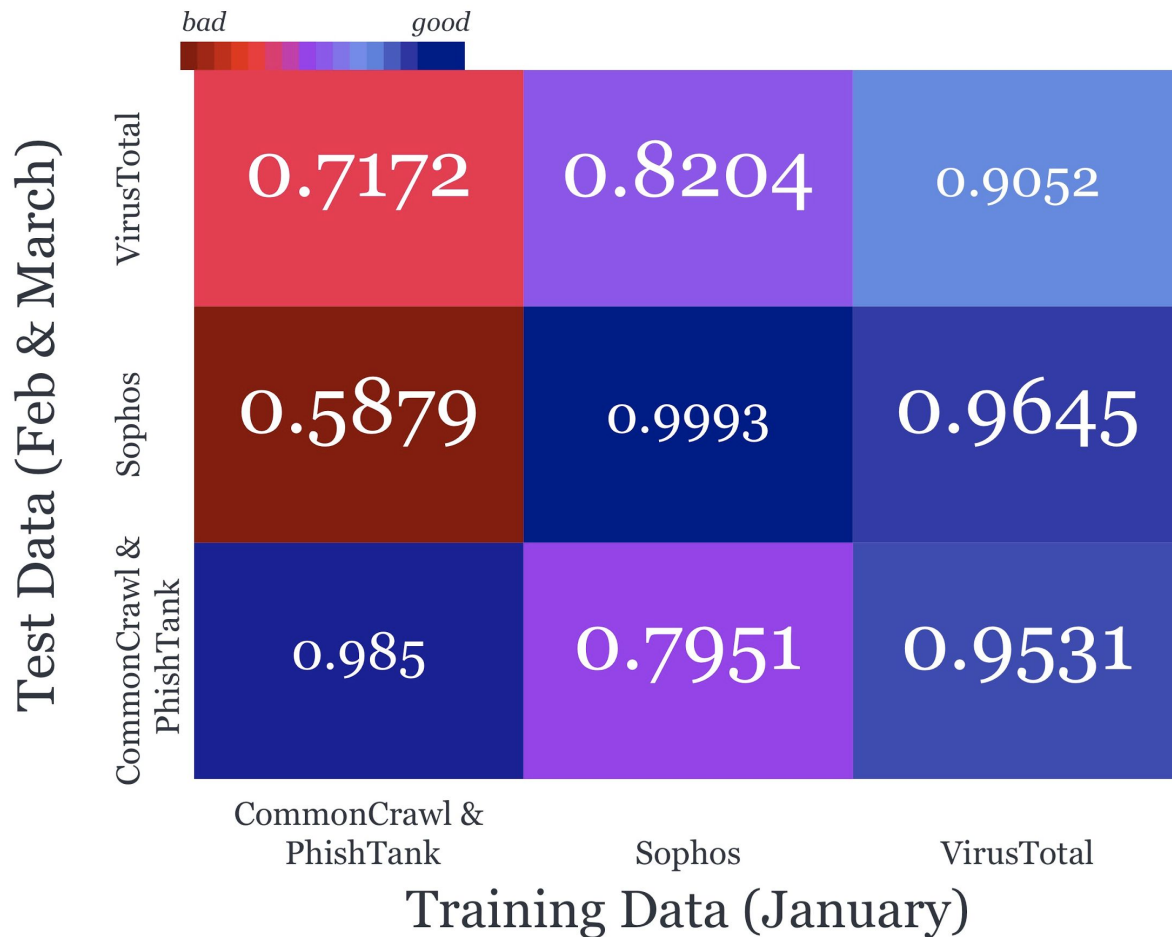2. Accuracy Metric Used: AUC

3. Datasets Used

4. Results!

AUC

Tested on VirusTotal

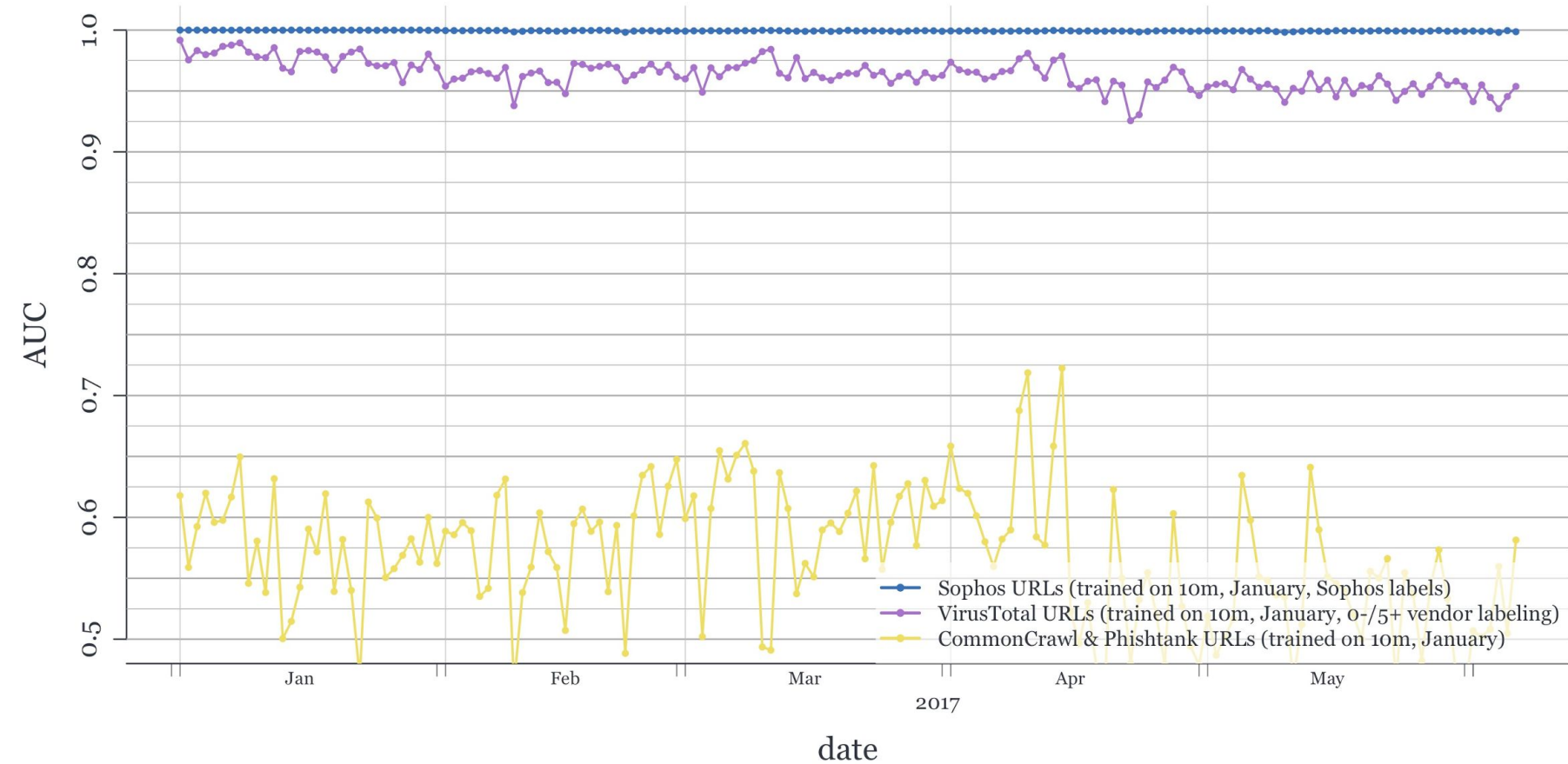Legend:
- Sophos URLs (trained on 10m, January, Sophos labels)
- VirusTotal URLs (trained on 10m, January, 0-/5+ vendor labeling)
- CommonCrawl & Phishtank URLs (trained on 10m, January)

**Tested on Sophos**

# What did we learn?

- **Model accuracy is *extremely* dependent on the training and test datasets used**

- Which datasets generalize better

- Expected variance in accuracy on new, inherently different data
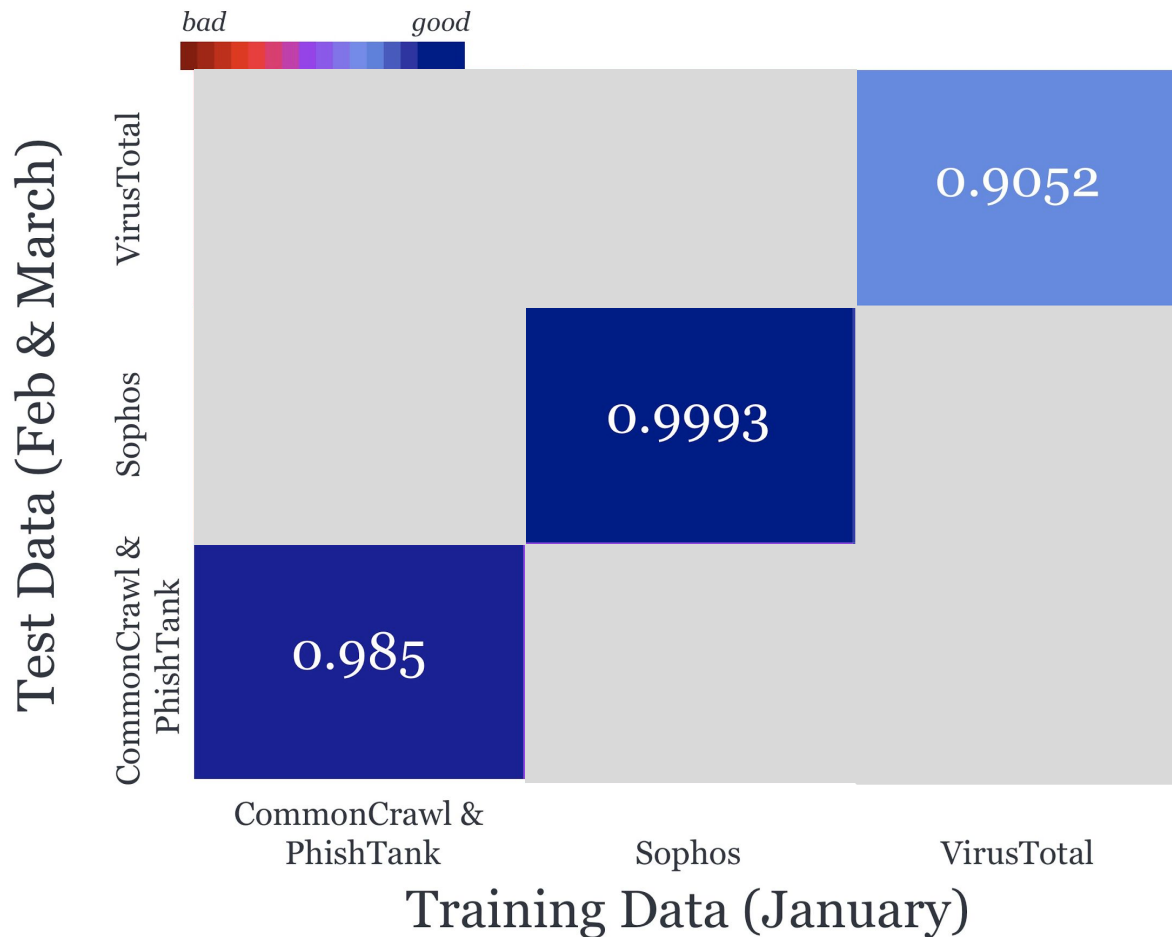
# What did we learn?

- Model accuracy is *extremely* dependent on the training and test datasets used

- **Which datasets generalize better**

- Expected variance in accuracy on new, inherently different data

# What did we learn?

- Model accuracy is *extremely* dependent on the training and test datasets used

- Which datasets generalize better

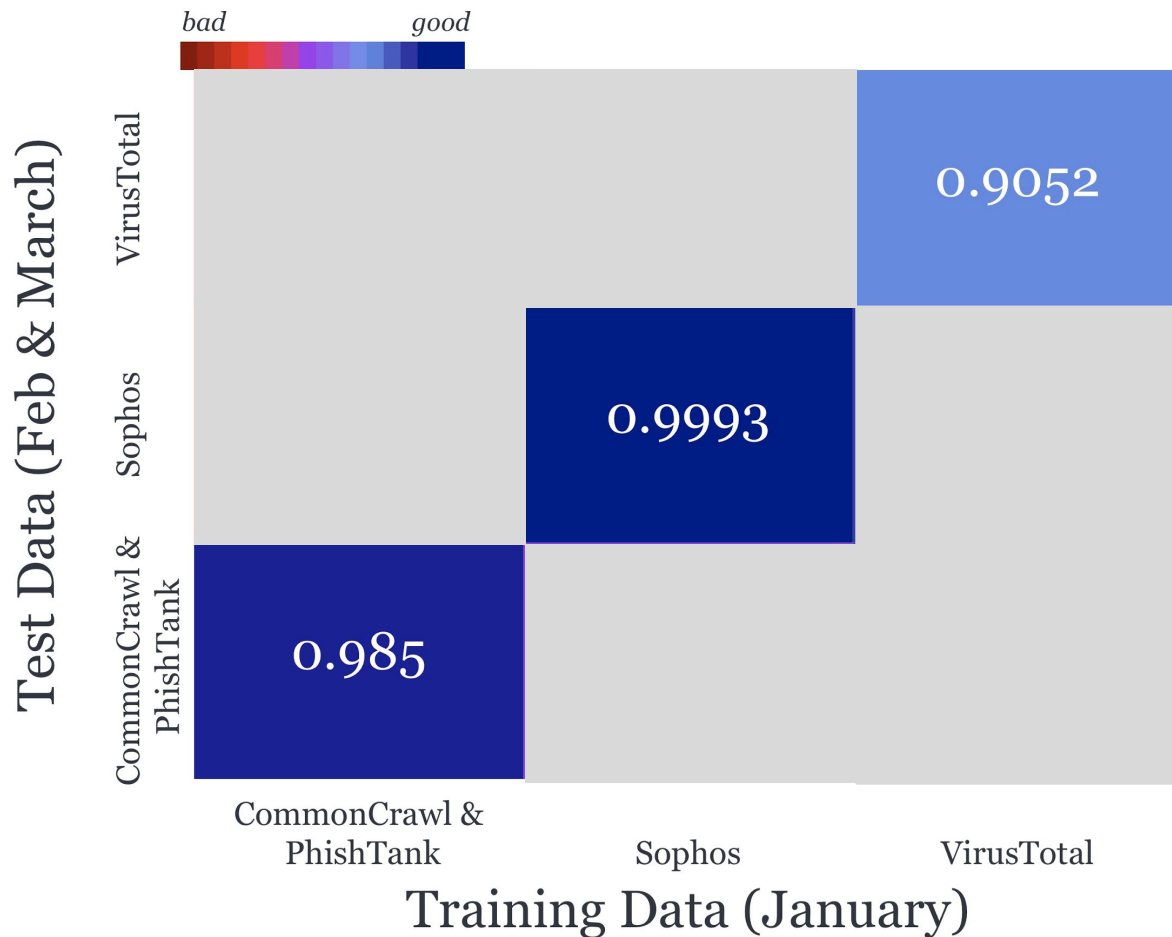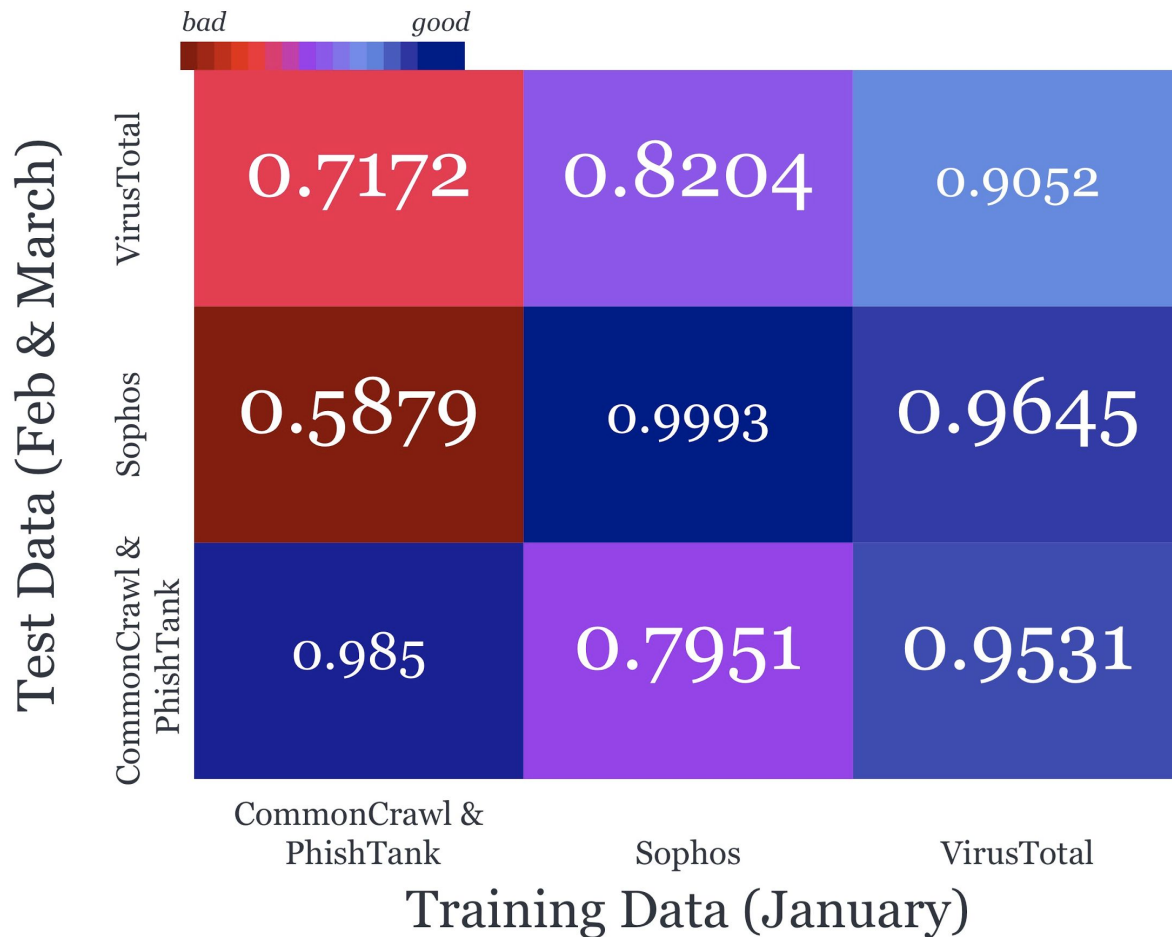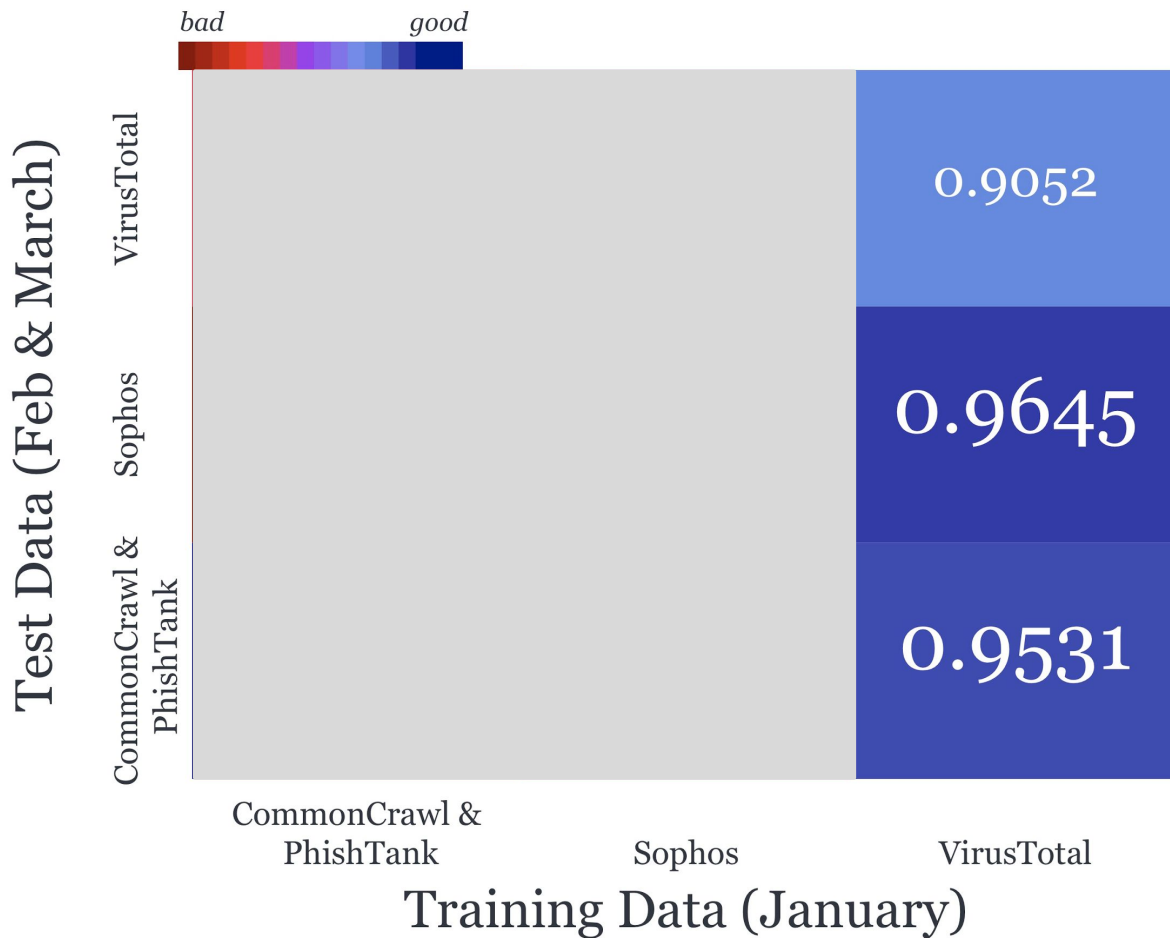- **Expected variance in accuracy on new, inherently different data**

# How minimize the probability… of failing spectacularly

## Models are liable to fail on different, future data.

*Especially when we lack deployment test data, we need to map the limitations of our models using **train / test dataset sensitivity analyses.***

***This technique can help us choose better training datasets and gain a better understanding of how sensitive model accuracy is to new test data distributions. This allows us to develop models that work in the <u>real world</u>, not just in idealized laboratory settings.***

# Thanks!

SOPHOS