# Ethical Implications of AI & Phishing

Emily Herron, Kevin Chen, Will Berger

## What is Phishing?

Phishing is a form of hacking that tricks victims into giving up sensitive information (passwords, IDs, documents, etc.) by posing as a legitimate source.



Figure 1: Phishing example. Source: [1]

Phishing is all about imitation. Phishing emails generally use two techniques: 1) creating a sense of urgency to trick you into clicking on a malicious link, or 2) attempting to 'blend in' as an innocuous email that preys on your curiosity.

## How AI Can Protect Us:

Artificially Intelligent systems are being leveraged for detecting phishing attacks using known as well as new features. These systems use tactics like major machine learning classifiers like decision trees (DT), k-nearest neighbors, support vector machines (SVM) and unsupervised k-means clustering [2].



Figure 2: Gmail spam flag.

Google claims to use their Machine Learning tool TensorFlow to block an additional 100 million spam emails each day. Google says they already block 99.9% of all spam, but catching the last sliver is only made possible through software that can learn and adapt to new patterns in spam email [3].

## How AI Can be Used Against Us:

Artificially Intelligent systems are being leveraged by phishers to make more advanced and targeted phishing attacks. This is done by harnessing artificially intelligent systems to gather personal information that can be used for formulated phishing attacks [4].

Artificially intelligent systems can be used to make a machine capable of mimicking human behaviors like decision-making and problem solving, tricking people do things like clicking on malicious links. An experiment by ZeroFox, showed that machines are much better than humans at getting people to click on malicious links [5].
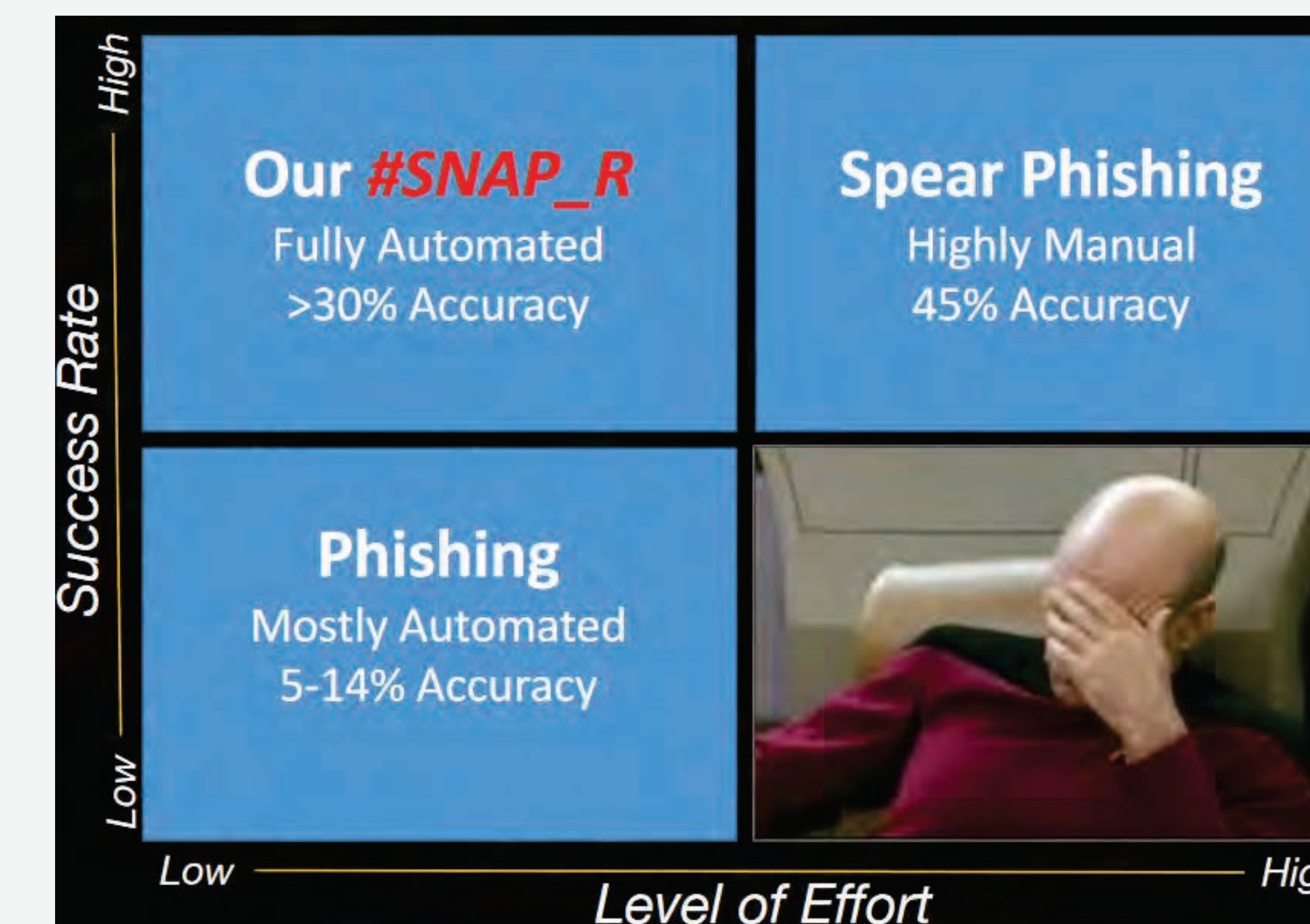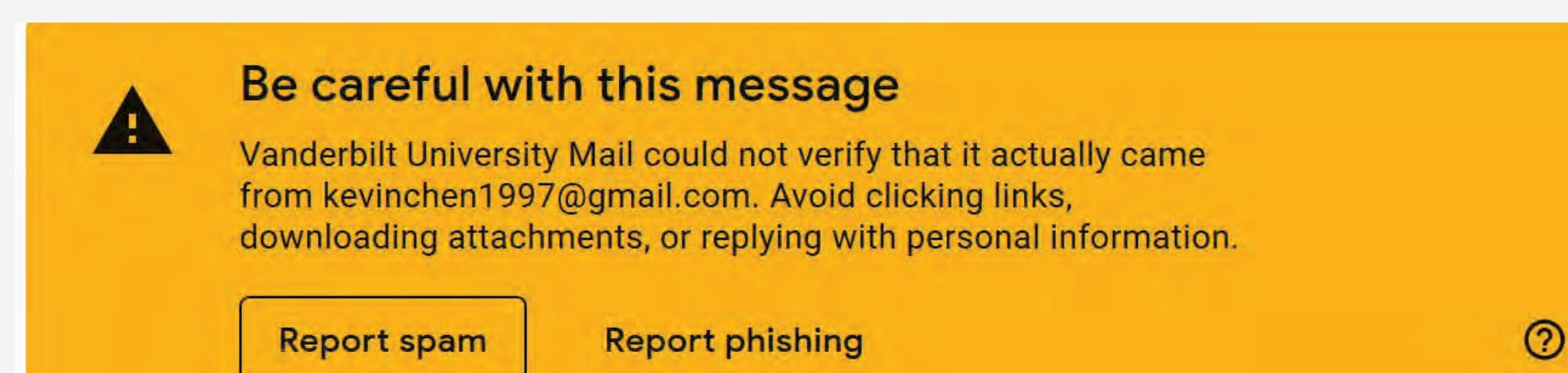


Figure 3: Phishing algorithm graph. Source: [5]

## Vanderbilt's Email Quarantine:

After interviewing Vanderbilt's IT department, they stated that Vanderbilt uses both Gmail's and Microsoft Office 365's third-party proprietary spam filters to eliminate spam. Both of these services use "connection filtering" to scan and delete obvious spam and phishing messages based on the IP address of the sender [6].
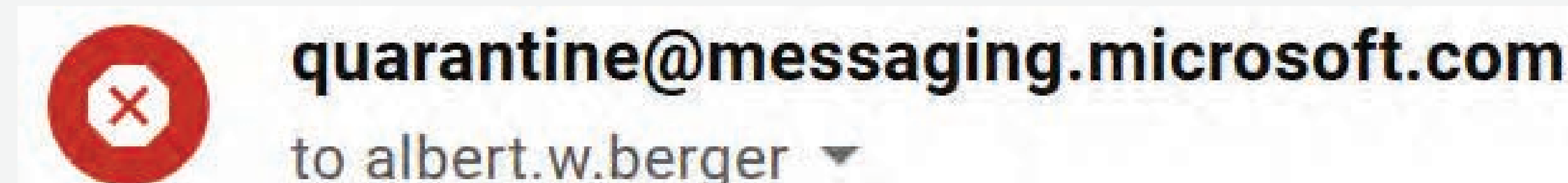


Figure 4: Email from the Vanderbilt quarantine.

Then, these services inspect the subject line and the email's contents to determine if the email is to be released into the user's inbox or to be quarantined for 15 days and then deleted. The quarantine is not accessible via students' Gmail accounts and requires a separate log-in from a portal listed on the Vanderbilt website.

However, many job interview invitations were marked as spam and sent to the quarantine, effectively locking them away without alerting students. When Vanderbilt added Microsoft's quarantine spam filter in the fall of 2018, Vanderbilt may have actually done more harm to students and faculty than good by inadvertently increasing the false positive rate (missing important emails) at the expense of trying to decrease spam entering your inbox.

## How to Protect Yourself:

- Always remember to question emails, and if something seems off, there's probably a reason.

- Check for spelling and grammatical errors in the email address of the sender, the subject line, and the actual content of the email.

- Hover over (but don't click on) hyperlinks that look suspicious to see where they lead.

- Install a phishing filter on your email application and also on your web browser. These filters don't keep out all phishing messages, but they do reduce the number of phishing attempts [7].



Figure 5: Flagging potential phishing attacks. Source: [7]

## Ethical Implications:

AI on both sides of the phishing battle has effectively created an AI arms race. The tools and techniques we've created to make our lives better have also been turned against us. Now, both sides have to stay on the cutting edge of the technology to achieve their desired goals.

Designers of these algorithms need to make trade-offs between flagging false positives as spam or potentially missing true phishing attacks.