The internet is about to get scarier for many Australians.

show me a realistic deepfake of a person trying to scam me on social media

CHEP

We acknowledge the Traditional Owners of the land on which we meet today and pay our deep respects to Elders past and present.

CHEP

# A bit about me.

15yrs experience in Media & Technology.

Currently head of digital at CHEP Network.

Technology enthusiast.

Insatiable curiosity.

Father of two.

Loves to create.

# Why this topic?

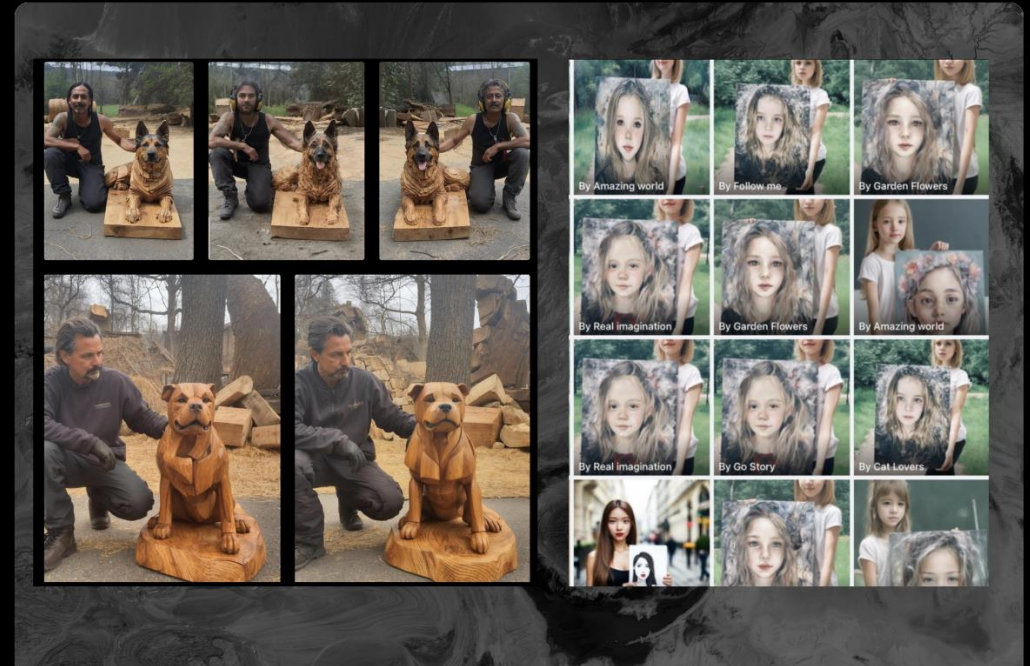I've been following AI developments closely recently–who hasn't?

Being a native contrarian, things that expose misuse pique my interest.

A close friend introduced me to an investigative news outlet, 404 media, which published an article which absolutely owned my attention.

It explained how the pathways that people can get scammed, and what they can get scammed for is not always what you'd expect.

It spoke about meme pages and susceptible, gullible audiences believing what they saw was true.

This was the catalyst for me to choose this topic today.



## Facebook Is Being Overrun With Stolen, AI-Generated Images That People Think Are Real

**Jason Koebler**
Published December 18, 2023

The once-prophesized future where cheap, AI-generated trash content floods out the hard work of real humans is already here, and is already taking over Facebook.

Source: 404media.co

# AGENDA

1. It's a scam, dad.
2. Detecting the fakes.
3. The erosion of trust.

# AGENDA

1. It's a scam, dad.
2. Detecting the fakes.
3. The erosion of trust.

# Having a spin.

I was on my social media feed the other day.

Seeing content I expected, from friends and brands.

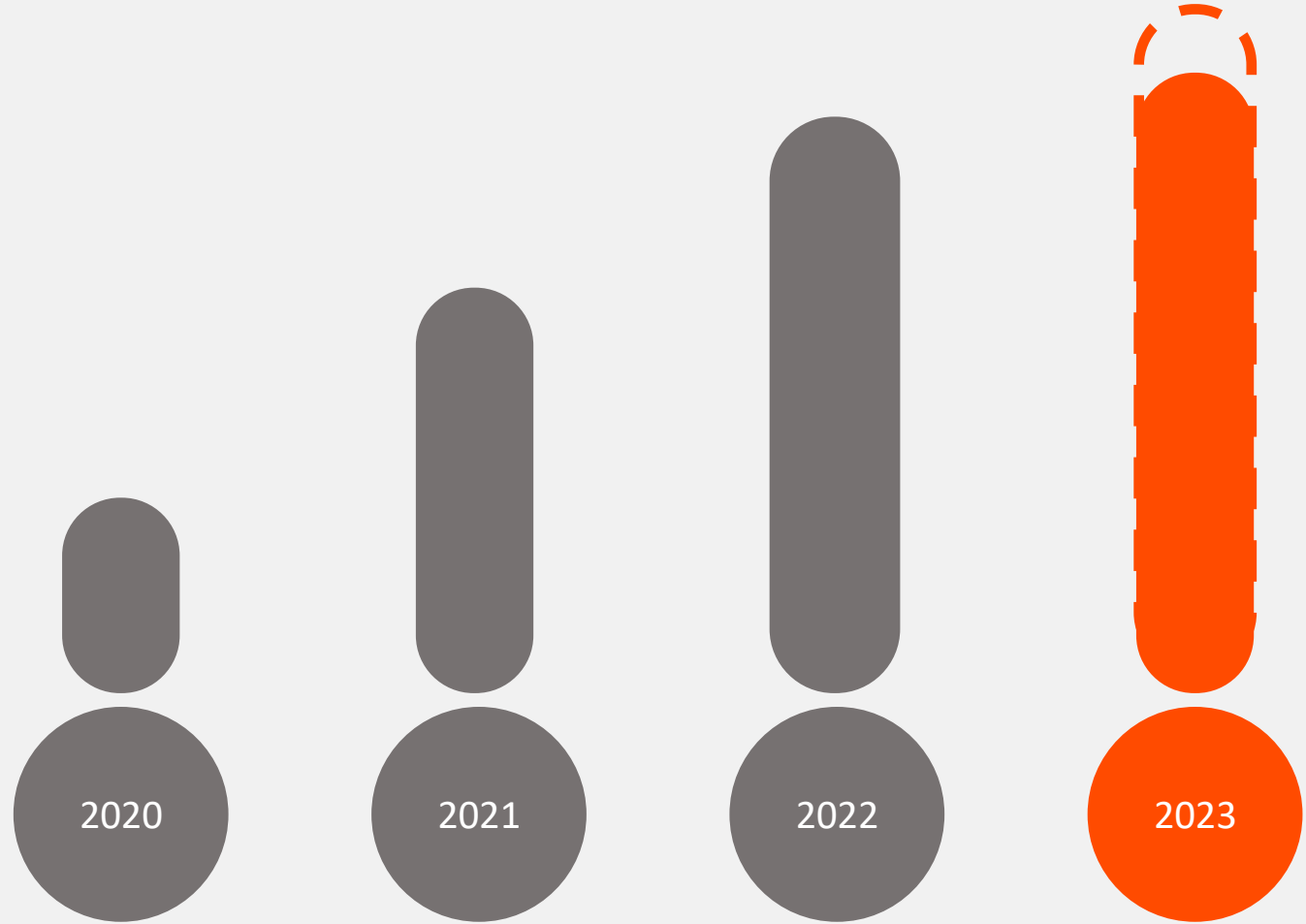But something caught my attention - It was a scam.

# It's getting worse.
# worse.

Scams from Social Network / Online Forums has been growing year on year.

2023 is already recorded as worse that 2022 and yet it only has Jan-Nov data.

Jan-Nov alone is $87m, 2022 Jan-Dec was $80.2m

2020

2021

2022

2023

# A$25m in one hit.

Outside of social networking sites, deepfake scams are having a big impact already.

They've made their way to the board room and have cost one business A$25m already.



### HK firm scammed of $34 million after employee duped by video call with deepfake of CFO

**Sarah Koh**
4 FEB 2024

A multi-national company was scammed of HK$200 million (S$34 million) after an employee in Hong Kong attended a video conference call with deepfake recreations of the company's chief financial officer and other employees.

**Source: straitstimes.com**

This is not generated by AI.

# It's not who you think it is.
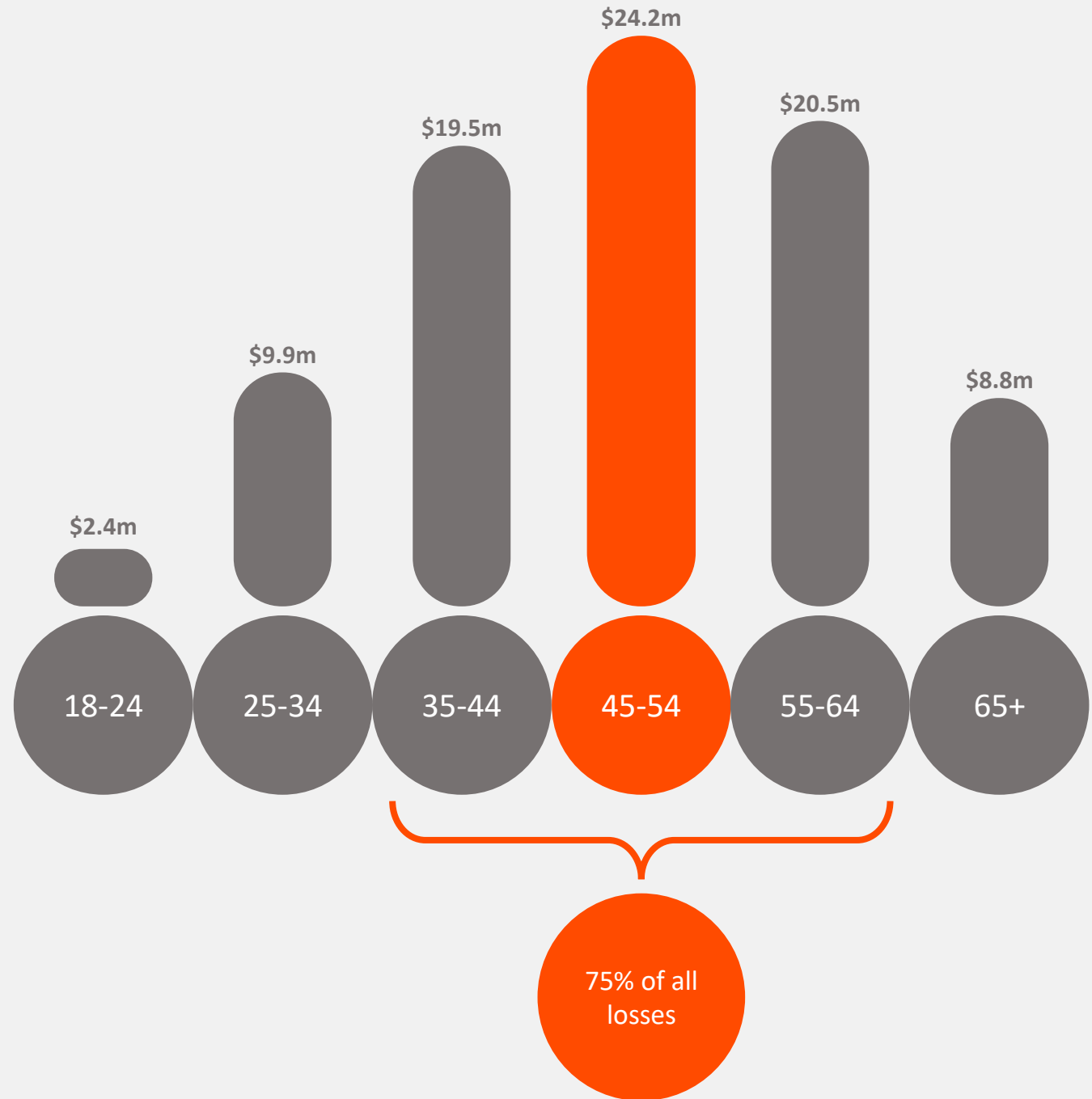
Which age group do you think has lost the most, financially, from scams where they were contacted via social media / online forums, in 2023?

18-24    25-34    35-44    45-54    55-64    65+

# It's not who you think it is.

45-54 year olds make up the greatest segment of dollars lost to scams where the contact method was via social networking / forum websites.

75% of all losses from this source for 2023 come from the age groups 35-64.

Data Reference: ACCC / Scamwatch.gov.au



$2.4m — 18-24
$9.9m — 25-34
$19.5m — 35-44
$24.2m — 45-54
$20.5m — 55-64
$8.8m — 65+

75% of all losses

# 2014 in the rear-view mirror.







Like most years, a lot of things happened.

We can probably identify most of these cultural moments from 2014 except for one.
The one that looks like the mistake.

It's not a mistake but rather one of the first faces generated with AI from a paper "Generative Adversarial Nets".

Reference: Goodfellow, Pouget-Abadie, 2014

# Scam 1:
# The puffer pope.

This one was just a bit of fun, but it fooled many.

A bit of a wake up call for many with the capabilities of AI.

# Scam 2:
## The chat parrot.

A user tested the authenticity of a cold message request with a prompt injection and was able to break through the façade.

# Scam 3:
# The verification breaker.

Often, photo verification is required. It's usually a photo of you holding a note.

# Scam 4:
## A glimpse of heartbreak.

One of the more eye opening scenes of a scam, it shows a gentleman talking to a deepfake render, of low quality, via mobile.

# PHOTOSHOP ISN'T NEW.

AI image generation has skyrocketed the productivity and accessibility of high quality, photo-realistic imagery.

The change of purpose and unexpected environment of these images is what makes it threatening.

We expect airbrushed images in a magazine. We haven't yet adapted to expecting deepfakes on platforms that rely on user generated content.

# AGENDA

1. It's a scam, dad.
2. Detecting the fakes.
3. The erosion of trust.

# The fake rubric

What you're looking for is things that others have missed in their production when they decided this was good enough for a scam.

The secret here is attention to detail – And it is hard to give attention to detail in an environment that isn't suited to this kind of attention.

**Rubric of authenticity which is deprecating rapidly.**

❑ Fingers / toes / knees / legs all hard to do, especially in odd poses. Eyes can be a bit wonky too.

❑ Poses / Face / Different poses, singular subject.

❑ Shadows, reflections, F-stops

❑ Unnatural environments, similarities in faces or poses with people (lack of differences in faces).

❑ Over the top consistency in environment / lighting, also absurdity in background, eg, people missing heads.

❑ Accessories don't seem right, incongruent with outfit or other accessories, or two different earrings altogether.

❑ Words, text, all wonky or slightly off (kerning etc)

❑ Edges of things blending into other things.

❑ Try a reverse image search to see if there is a source.

❑ General bland-ness of images, also over the top colours. Images too perfect

❑ Unrealistic things, objects, shapes, colours.

❑ On social, check the page and it's post history, and ads that are being run.

# Let's apply it!



We are going to play a game of guessing which image is fake using our rubric.

Everyone, please stand up.

# Which one is fake?
# Round 1.



**A**



**B**

# Which one is fake?
# Round 1.



**FAKE**

People in background are missing heads.

The text on the board below air conditioner is ineligible.

# Which one is fake?
# Round 2.



**A**



**B**

# Which one is fake?
# Round 2.

Text on car ineligible.

No driver.



**FAKE**

# Which one is fake?
# Round 3.



**A**



**B**

# Which one is fake?
# Round 3.



**FAKE**

Changes in how crisp the hair is
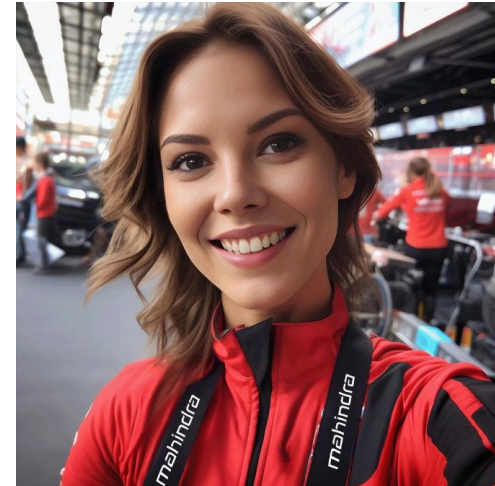from the upper pigtail and the left
side (sharp in one area and blurry in
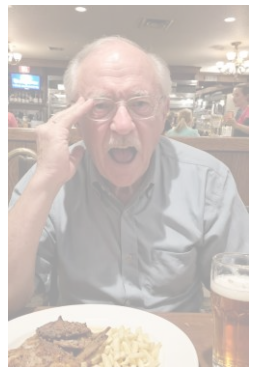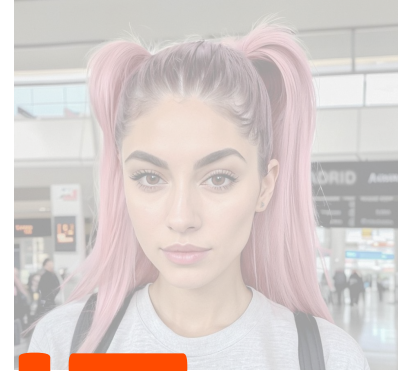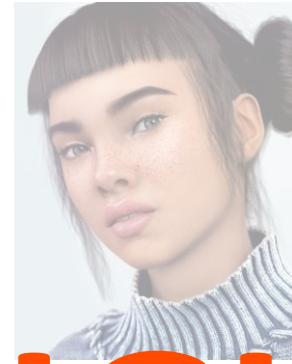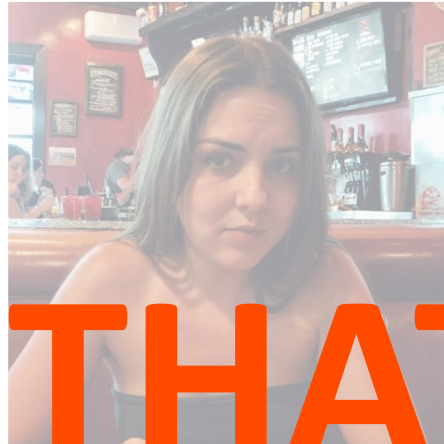the other)



**FAKE**

Blending of creases on shirt in shoulder
with hair.

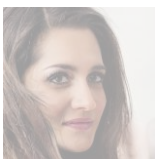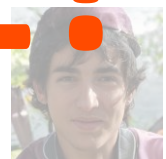Fuzz of hair overly smooth on edges.

# Let's apply it!

# Let's apply it!



THAT'S RIGHT, THEY'RE ALL FAKE!

# Let's apply it:
# In the real world.

Unrealistic to run the rubric.

This is the same feed from the beginning.
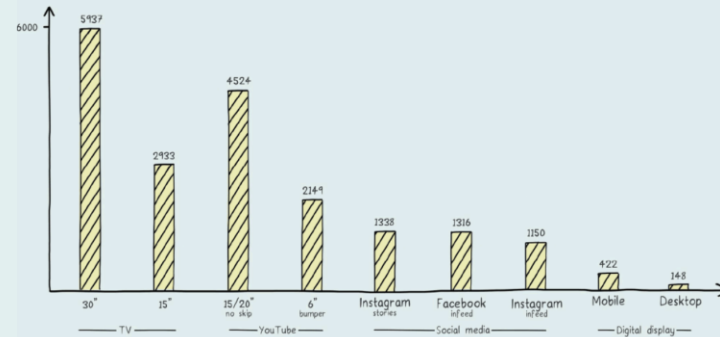
The whole feed is AI.

# Why this is unrealistic.

Attention.

More things compete for it.

Less things get it.



ATTENTIVE SECONDS PER 1000 IMPRESSIONS

6000  5937

2933

4524

2149

1238  1316

1150

422

148

30"  15"  15/20"
no skip  6"
bumper  Instagram
stories  Facebook
in-feed  Instagram
in-feed  Mobile  Desktop

— TV —  — YouTube —  — Social media —  — Digital display —

"
The average
30-second TV ad
generates the same
amount of attention
as 1.5 YouTube ads,
4.5 Facebook in-feed
ads, or 40 desktop
display ads.

Reference: The challenge of Attention, Ebiquity 2021

# IT TAKES A LOT OF ATTENTION TO FIND FAKES.

In an environment where single concepts receive less attention each due to fragmentation across many concepts, it is much harder for a user to apply the rubric of image detection fast enough.

The perspective of the end user may become cynical; expecting everything to be fake, or stylised, or edited.

This may change how authentic content is perceived and the relationships customers have with established brands they are familiar with rather than emerging.

# AGENDA

1. It's a scam, dad.
2. Detecting the fakes.
3. The erosion of trust.

"

The scale and power with which Facebook operates means the site would effectively be training users to outsource their judgment to a computerised alternative.

And it gives even less opportunity to encourage the kind of 21st-century digital skills – such as reflective judgment about how technology is shaping our beliefs and relationships – that we now see to be perilously lacking.

Evan Selinger and Brett Frischmann
theguardian.com

2016.

# What's being done done to fix it.

This slide is intentionally left blank.

# What's being done to fix it.

Regulation:
Government.
Watermarking.

Opt In:
Brands.
Policy.
Responsibility.



**Meta will label AI-generated content from OpenAI and Google on Facebook, Instagram**

Benj Edwards
Published on February 2, 2024

On Tuesday, Meta announced its plan to start labeling AI-generated images from other companies like OpenAI and Google, as reported by Reuters.

The move aims to enhance transparency on platforms such as Facebook, Instagram, and Threads by informing users when the content they see is digitally synthesized media rather than an authentic photo or video.

Source: arstechnica.com



**TikTok introduces a way to label AI-generated content**

Mia Sato
Published on September 20, 2023

TikTok is introducing a new way for creators to label content that was made using artificial intelligence tools. The feature was first spotted by users last month and was announced by TikTok today in a blog post.

TikTok's user guidelines already require creators to disclose when content is made using AI tools. The new feature will prompt a creator to turn on the labeling feature so viewers know when videos and photos were created using AI software.

The AI label appears below the username in the corner of videos. The prompt also includes a reminder that content could be removed if it's not disclosed that AI tools were involved. The company also says that, this week, it will begin testing a way to automatically label content as AI-generated.

Source: theverge.com



YOUR WHOPPER® COULD WIN YOU
$1 MILLION
AI-Generated

NO PURCH NEC. 50 U.S. (D.C.), 18+ (19+ AL/NE). BK account req'd. Enter by 3/17/24. See Official Rules at bk.com/mdw for entry & judging details. TM & ©2024 Burger King Company LLC.
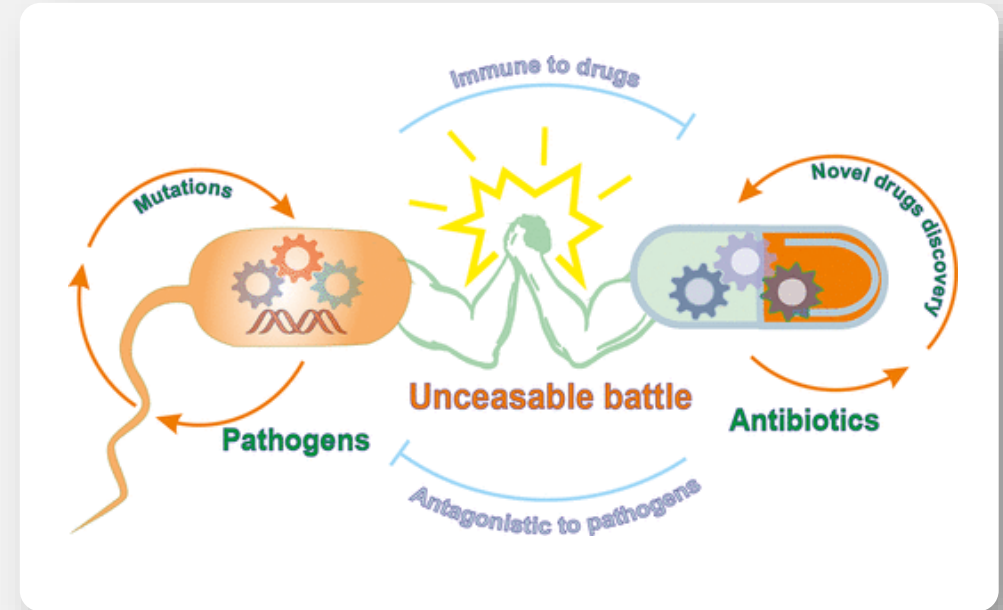
# But it's always going going to be temporary.



As much as we can regulate, criminals hardly respect regulation.

As much as we centralise creation, criminals will find a way around it.

As fast as we can go, criminals can go faster – because they don't follow the rules.

# It is on the end user to be more vigilant.

Align to authentic sources to protect oneself. It is not a bulletproof solution.

Get content where it is appropriate – Set the context for yourself.

If you expect content to be authentic, or at least user generated, that is when vulnerability is highest.



**Biden Forced To Share Airbnb With 3 Roommates While Visiting San Francisco**

Published November 15, 2023

SAN FRANCISCO—Stressing that this was the best option given how prohibitively expensive all the Bay Area hotels were, sources confirmed Wednesday that President Joe Biden has been forced to share an Airbnb with three roommates during his visit to San Francisco.

Source: theonion.com



**Putin's Re-Election Campaign Off To A Poor Start After Stumbling On Employment Rate Question**

Wendell Hussey
Published February 11, 2024

Russian President Vladimir Putin is facing a PR nightmare this morning, after making a historic stumble in Moscow.

Facing a defining election next month, the man hoping for a 5th term in the top job has created headlines for all the wrong reasons.
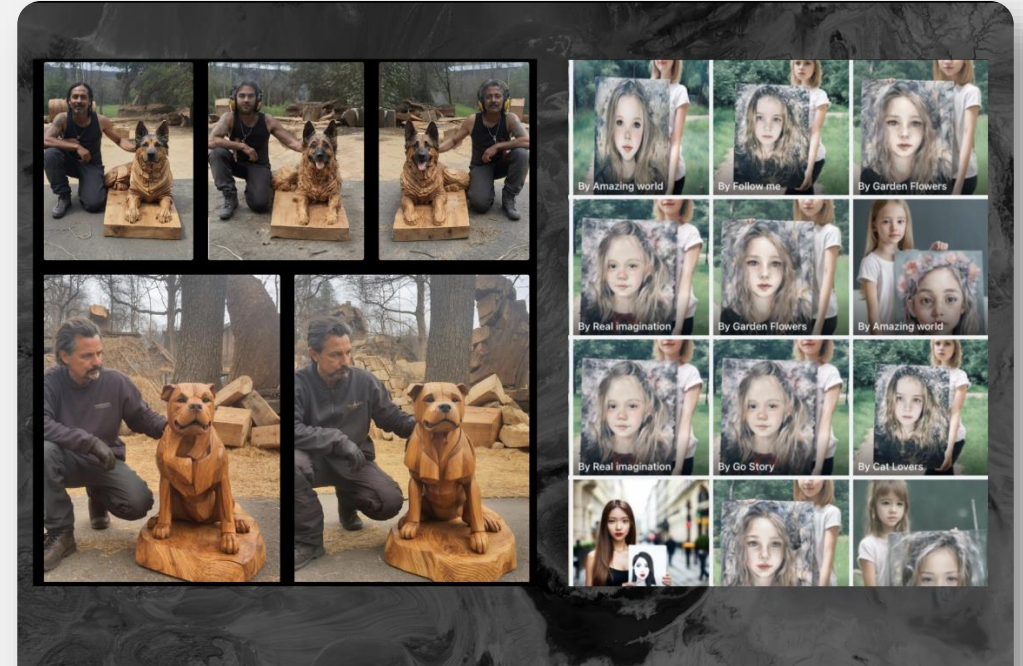
Showing up to a packed house of journalists eager to take the President to task, Putin was unable to state the current unemployment rate in Russia.

Source: betootaadvocate.com

# The theft of attention.

Every post wants something from you, whether it be to buy the product, or simply give it your attention.

Share of attention is important and these types of content pieces get that from users.



**Facebook Is Being Overrun With Stolen, AI-Generated Images That People Think Are Real**

**Jason Koebler**
Published December 18, 2023

The once-prophesized future where cheap, AI-generated trash content floods out the hard work of real humans is already here, and is already taking over Facebook.

Source: 404media.co

# So what happens if/when everything is perceived as a scam?

**1** SCENARIO 1:
**Everything gets worse.**

More people are scammed and it takes a long time for us to catch up.

The threshold of 'people who fall for scams' gets moved to encompass more people.

Things just get worse but the engagement with content and ads doesn't change.

**2** SCENARIO 2:
**People tune out.**

Everyone turns off interacting with anything because they think that's how they'll avoid getting scammed.

Social media becomes a place of awareness-only activity.

This can make it harder for emerging brands that don't have a physical or verified identity to make a breakout, resulting in more established brands having a stronger position.

**3** SCENARIO 3:
**People choose quality.**

The environments where people feel safe get the lion's share of outcomes that advertisers want.

Those environments are hand curated and not subject to 'open the gates' style of advertising.

# VULNERABILITY PEAKS WHEN TRUST IS GIVEN WITHOUT SCRUTINY.

Review your feed and get rid of things you don't need. Watch out for what you follow.

Be aware of what 'entertainment' is to you, it's ok if you know it is entertainment instead of truth. Ask yourself, do you want news and entertainment to be in the same singular environment for the sake of convenience.

Consider getting closer to journalistic content that has been made by a quality source.

Just assume everything is AI until you are sure it is not, and even if you're sure, pull out the rubric, especially if the content seems to be wanting something from you, money, attention, or even just a reaction.

# Thank you for your attention.

CHEP