



SPRITZ  
SECURITY & PRIVACY  
RESEARCH GROUP

# User Profiling in Video Games: From Identification to Private Data Inference

**Pier Paolo Tricomi**

*tricomi@math.unipd.it*

University of Padua, Italy

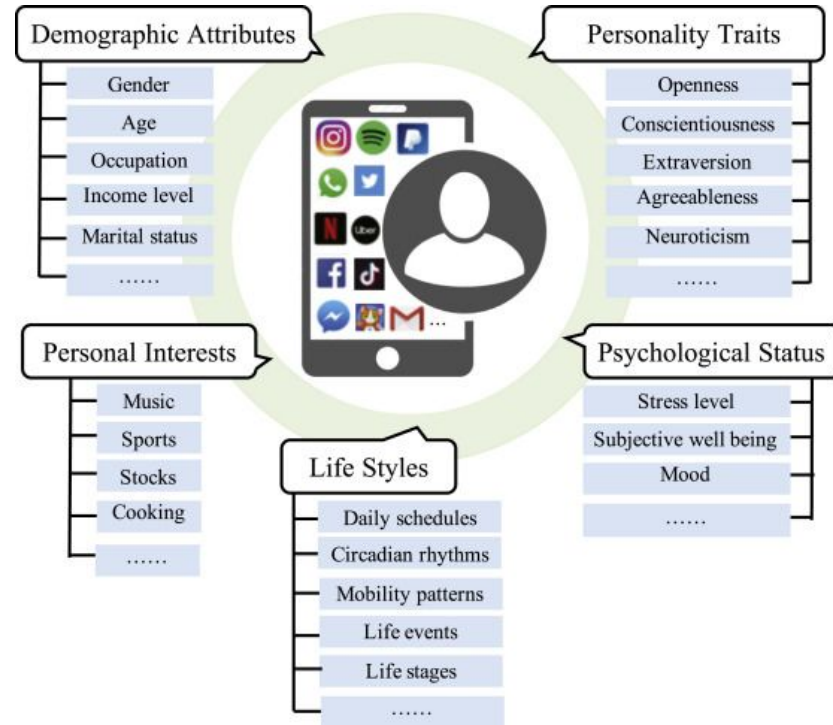
18<sup>th</sup> October 2023



# User Profiling?

Create a profile of your users, for:

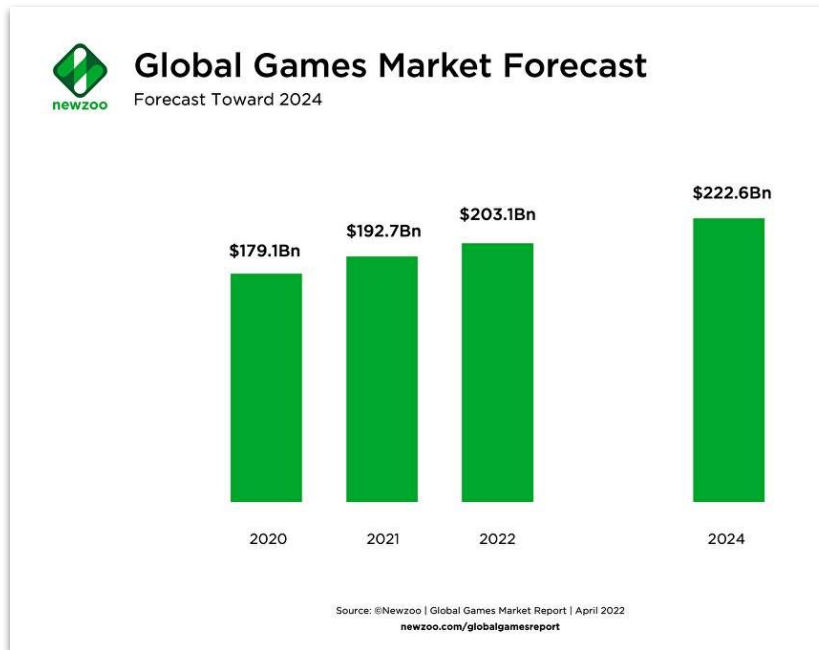
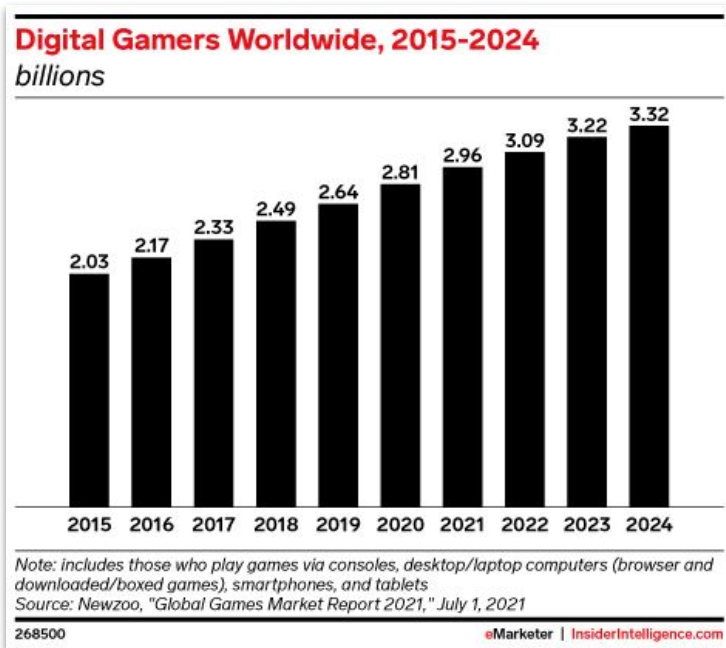
- Customers analytics
- Marketing strategies
- Custom experience
- Sell data (!)
- ...



# The problem (1)

3.09 billion people on Earth are gamers (2022)

Video game market generated 200\$ billions (2022) (Amazon + Meta + Google)



# The problem (2)

A lot of money involved -> Scams, account take over

Many people involved -> Profiling, malicious activities

**CBC** | MENU

news Top Stories Local The National Opinion World Canada

Technology & Science

## PlayStation data breach deemed in 'top 5 ever'

Emily Chung - CBC News - Posted: Apr 27, 2011 10:59 AM ET | Last Updated: April 28, 2011

Names, birthdates and some credit card data may have been stolen from users of Sony's PlayStation Network in what may be one of the biggest data breaches ever.

More than 75 million accounts worldwide, including more than one million in Canada, are registered with the network that suffered a massive data breach this week, Sony confirmed Wednesday.

The massive breach is one of the "top five ever," said Alan Paller, director of research for the SANS Institute, a cybersecurity training and research institution based in Bethesda, Md.

### Scam of the Week – Millions of League of Legends Gamers Targeted with Phishing Scam

November 1, 2018 10:13 am Geraldine Strawbridge



Hundred million players worldwide, the game has been voted the most popular video game and has attracted the attention of cybercriminals who are keen to take advantage of this

people as they can fall for their scam, the criminals have sent users an email with a phishing login page which is almost identical to the real site. The branding, layout and design have been replicated, however the site is just a fake phishing website set up to harvest user

The vast majority of people will tend to use the same login details and passwords for many accounts, so if they manage to obtain the details for one, they have free reign to break into

07/05/2021

## CONCLUSION PAPER

RAN C&N – Digital grooming tactics on video gaming (adjacent) platforms  
15-16 March 2021, Online event

## Digital Grooming Tactics on Video Gaming & Video Gaming Adjacent Platforms: Threats and Opportunities

**VB**

f t i in F

## Cyber-bullying and video games

Jesse Aaron@JesseAarone September 26, 2014 6:56 PM

The saying goes, "kids will be kids." But what does the life of a modern kid in American, European, and Asian countries share?

Internet access.

Webpage: ec.europa.eu/ran

twitter | facebook | linkedin | youtube

physical, school-ground kind – has actually been classified as a social issue. It wasn't until the mid-90s where instances of cyber-bullying drew enough attention to warrant a subset classification, primarily because cyber-bullying were either identical or more severe than physical bullying. Cyber-bullying can include emotional stress, self-harm, and in rare cases, murder or

that socially conscious parents and politicians have been pushing for anti-bullying laws, it appears as though the anonymity that the internet is giving bullies another outlet to harass their prey. It's easy for bullies to flock to cyberspace and online video games to bully. It's convenient for them. The prevalence is so great that the sheer

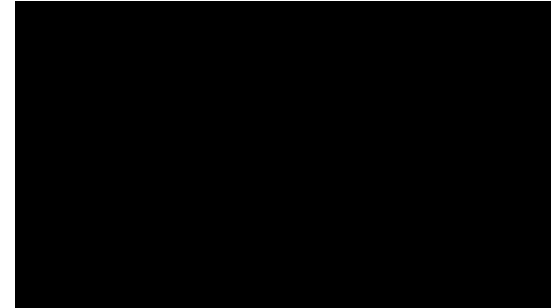
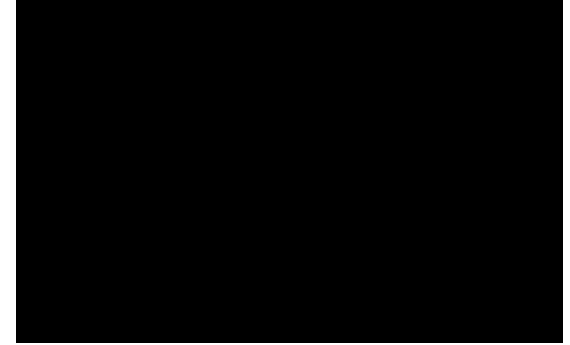
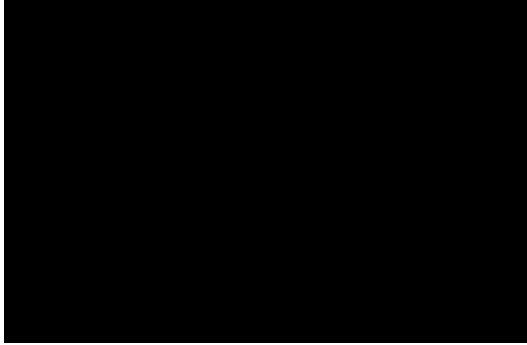
# What to do?

**Problems can be reduced uniquely Recognizing/Identifying a player:**

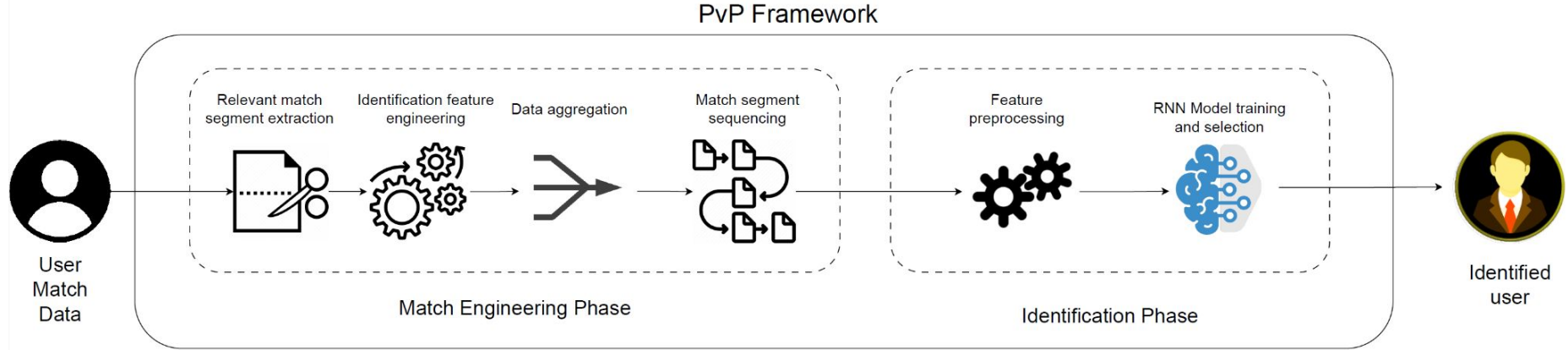
- Create a game “fingerprint”
- Ban harmful players from all their account
- Create new “biometric” authentication system

**The fingerprint can be the gamer play-style!**

# Identification intuition - Movement Action & Camera



# The Identification Framework



# Player Identification – Dota 2 Dataset Creation

- Survey to collect players and their data (matches replays)
- 50 players, 100 matches per player, and 5000 matches in total (balanced dataset)
- Sequences of states (cursor and camera positions) and actions (attack, move...)
- Sequences of 2 minutes

Type	Features
Cursor, Camera Cell	X_mean, X_std, X_changes Y_mean, Y_std, Y_changes
Camera Vector	X_mean, X_std, X_changes Y_mean, Y_std, Y_changes Z_mean, Z_std, Z_changes
Action: Move_to_position	n_occurs, X_mean, X_std, Y_mean, Y_std
Action: Move_to_target	n_occurs
Action: Attack_move	n_occurs
Action: Attack_target	n_occurs
Action: Cast_position	n_occurs
Action: Cast_target	n_occurs
Action: Cast_target_tree	n_occurs
Action: Cast_no_target	n_occurs
Action: Hold_position	n_occurs
Action: Drop_item	n_occurs
Action: Ping_ability	n_occurs
Action: Continue	n_occurs

Features



# Player Identification – Preliminary Model

## Preliminary Model:

- Two LSTM layers (64 unit each, *tanh*)
- Fully connected layer (64 unit, ReLU)
- Output: Softmax layer (50 unit)

Categorical Crossentropy Loss function

Adam optimizer (learning\_rate = 0.001)

Batch\_size = 256, 100 epochs

Model: "sequential"

Layer (type)	Output Shape	Param #
lstm (LSTM)	(None, 240, 64)	26112
lstm_1 (LSTM)	(None, 64)	33024
dense (Dense)	(None, 64)	4160
dense_1 (Dense)	(None, 50)	3250

Total params: 66,546

Trainable params: 66,546

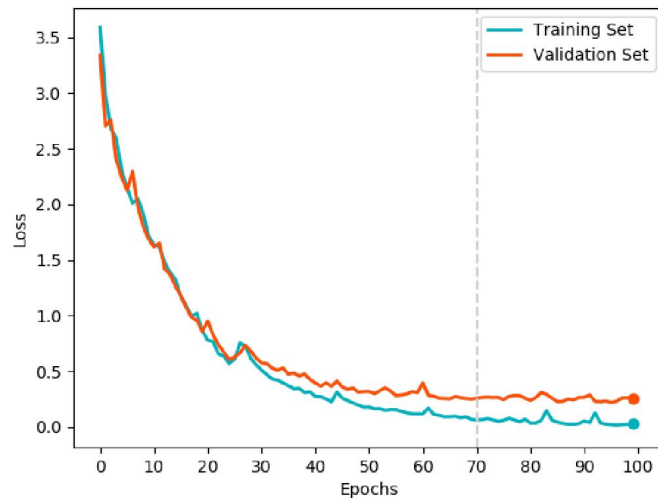
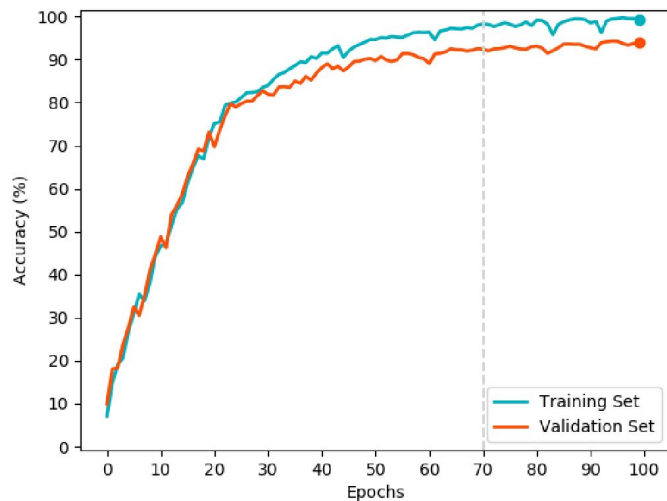
Non-trainable params: 0

Preliminary Model Summary

## Player Identification – Model selection (2)

Good generalization, low risk of overfitting

Stabilize after ~ 70 epochs



# Player Identification – Evaluations

**Best model on Validation Set** (Accuracy = 96.48%, loss = 0.179):

- 256 units both LSTM layers
- 128 units dense layer
- learning rate = 0.001

On **test set: accuracy = 96.32%, loss = 0.198**

Very high generalization, play-style can be considered “unique”

Using only cursor, camera and move action (**common features**):  
Accuracy = 95.6%, Loss = 0.162

# Player Identification – CS: GO Case Study

**50 Players, First 10 minutes, 100 matches each**

**On test set:**

- **All Features 91.68% Accuracy**
- **General Features 85.83% Accuracy**

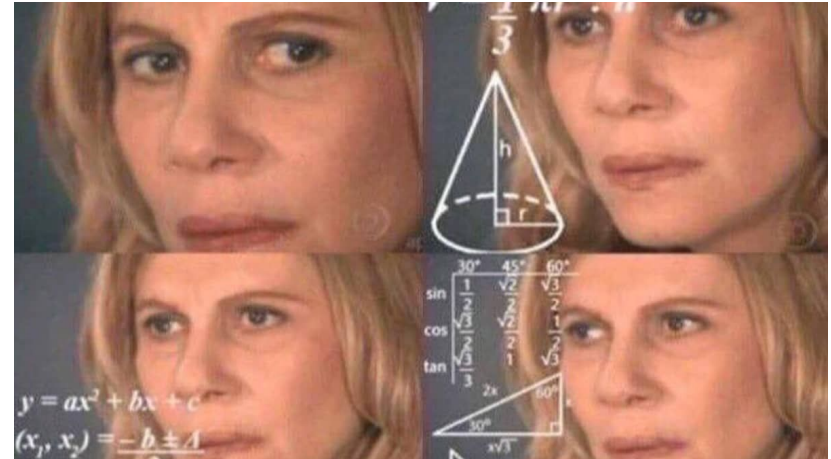
Framework generalizes well!

# Recap

Identification is possible...

Data is publicly available...

Mmmhh...



Can we infer more information about a video gamer?

# Using Machine Learning to violate the Privacy of Video Gamers

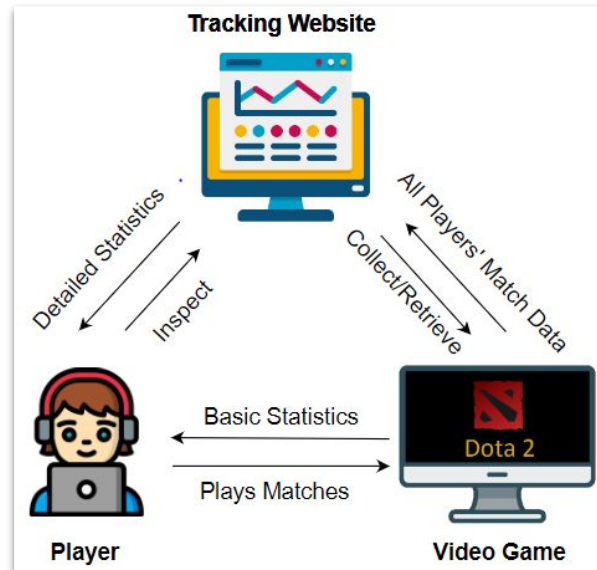
Pier Paolo Tricomi, Lisa Facciolo, Giovanni Apruzzese, Mauro Conti. “*Attribute Inference Attacks in Online Multiplayer Video Games: a Case Study on Dota2.*” CODASPY 2023

## More Context (1)

- Video Games (VG) are becoming increasingly popular
  - One of the few industries that are constantly improving their profits
- Some *competitive* VG are denoted as “E-sports”
  - Examples: Dota2, Fortnite, League of Legends
- Some tournaments of such E-sports have very high prize-pools
  - For Dota2, “The International” had a prize pool of 40M \$ in 2021

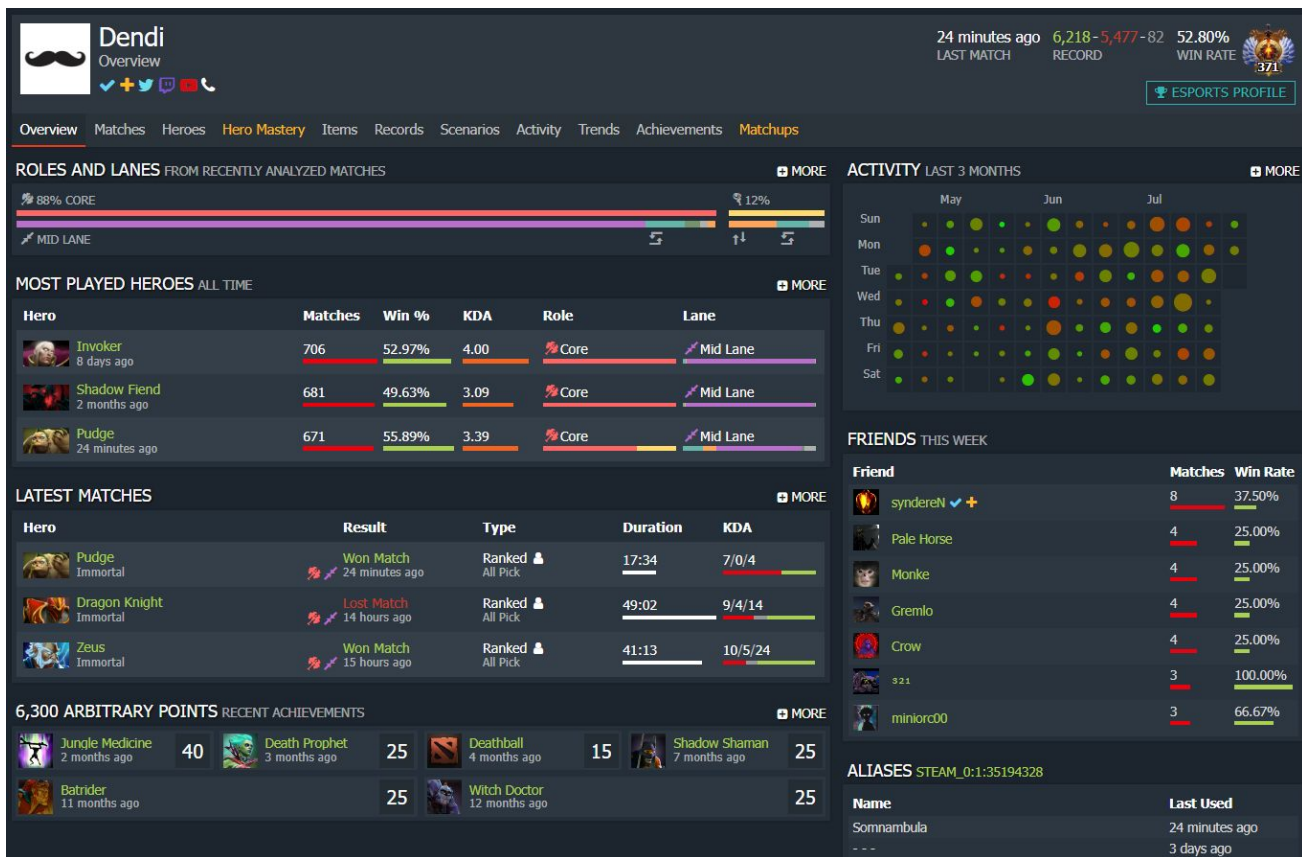
## More Context (2)

- Such prizes attract a lot of players who “play-to-win” and want to get better...
  - Best way of improving at something? Learn from past mistakes!
- ...which, in the E-sport ecosystem, it can be easily done via Tracking Websites





# A Tracking Website



# A Tracking Website



Dendi  
Overview

24 minutes ago  
LAST MATCH

6,218-5,477-82  
RECORD

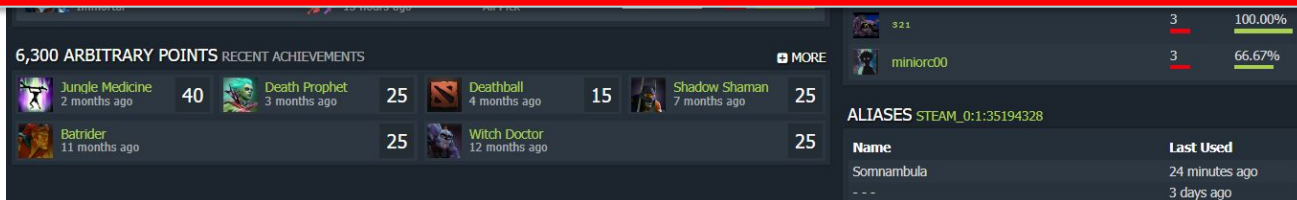
52.80%  
WIN RATE

321  
ESPORTS PROFILE

All of this is Public

—

for 70M DOTA2 players



6,300 ARBITRARY POINTS RECENT ACHIEVEMENTS MORE

Achievement	Points	Time Ago
Jungle Medicine	40	2 months ago
Death Prophet	25	3 months ago
Deathball	15	4 months ago
Shadow Shaman	25	7 months ago
Batrider	25	11 months ago
Witch Doctor	25	12 months ago

ALIASES STEAM\_0:1:35194328

Name	Last Used
Somnambula	24 minutes ago
---	3 days ago

# Why Public?

**It is the playerbase who want the statistics collected by TW to be publicly available!**

The reasons are various, e.g.:

1. Inspecting the profiles of *other* players can be used to learn some of their tricks...
2. ...in turn, by having their own profile publicly accessible, a given player can gain visibility if they perform very well...
3. ...such “visibility” can lead to invitations to play in top-teams, or to finding new (good) teammates
4. The visibility can come either because other players “inspect” a given player’s profile, or because of climbing “public ladders”

# All such data is public, OK... so what?

I don't have any problems if others know:

- that I win very often...
- ...or that I regularly play with a given hero...
- ...or that I adopt an aggressive playstyle...
- ...or that I communicate in the chat by using DOTA2 jargon...
- ...or that I frequently play on a given day of the week...

**...right?**

# All such data is public, OK... so what?

I don't have any problems if others know:

- that I win very often...
- ...or that I regularly play with a given hero...
- ...or that I adopt an aggressive playstyle...
- ...or that I communicate in the chat by using DOTA2 jargon...
- ...or that I frequently play on a given day of the week...

**Problem:** such “availability” exposes E-sports’ players to the risk of “Attribute Inference Attacks” (AIA)

# Attribute Inference Attack 101

Use Machine Learning to Infer Private Data from Public Data

## **Assumptions:**

- In a specific environment, everyone release some public data
- Some people release their “private” data publicly as well (e.g., age, gender)

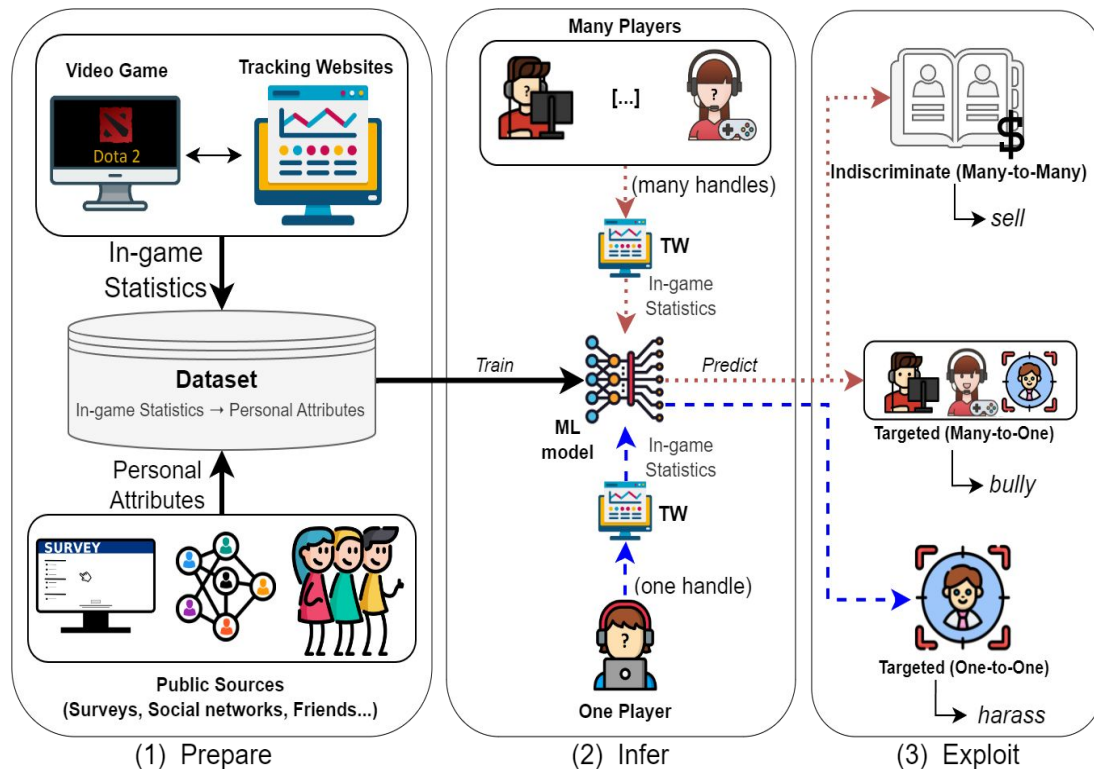
## **Method:**

- Train a machine learning model that maps public data to private data
- Use such model to predict private data of people not releasing it

## **Examples:**

- Social Media
- E-commerce/streaming platform ratings

# Our Threat Model



# Our Assessment (1)

- We proactively assess such a threat, because *nobody* ever did something similar in the E-Sport ecosystem. We focus on Dota2
- We conduct an informed survey, asking ~500 Dota2 players to provide us with private (non-sensitive) information about their real-life (e.g., age, gender, occupation, whether they buy Dota2 content, and OCEAN personality traits)
- We use the handle (i.e., nickname) of such players to collect their (publicly available) Dota2 in-game statistics from popular TW (opendota).



# Our Assessment (2)

- We **find a correlation (!)** between the players in-game statistics and their real life.
  - Such a finding suggests that AIA can be successful!

Gaming and Private Information correlation was already proved in other genres (mainly RPG)



## Our Assessment (3)

- We (ethically) perform diverse AIA: we use 80% of our data to train ML models, and predict the personal attributes of the players included in the remaining 20%.
  - Player Level (P): Consider aggregated statistics of one month
  - Match Level (M): Consider all **single** matches in the last month
  - Reduced Match Level ( $M^-$ ): consider at most 30 **single** matches in the last month
- Why need M &  $M^-$  ?

## Our Assessment (3)

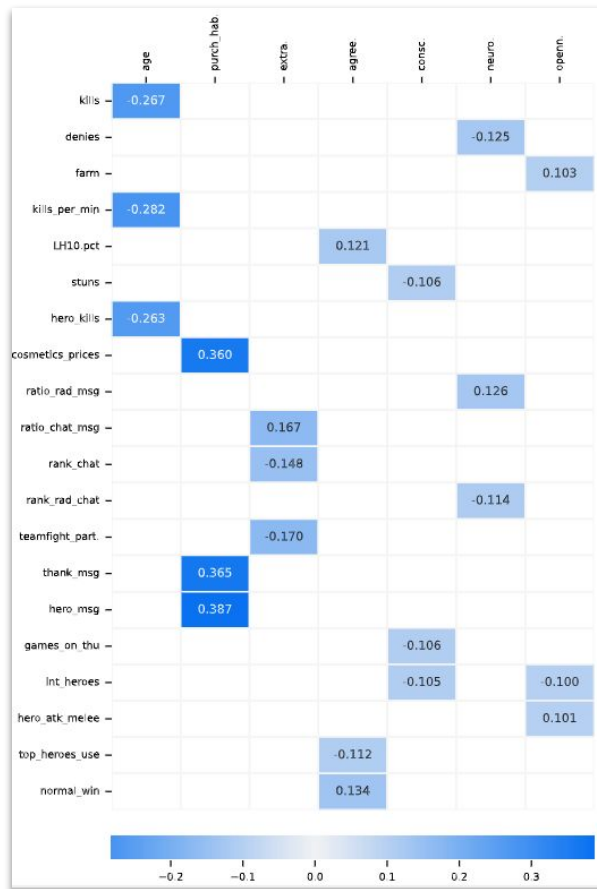
- We (ethically) perform diverse AIA: we use 80% of our data to train ML models, and predict the personal attributes of the players included in the remaining 20%.
  - Player Level (P): Consider aggregated statistics of one month
  - Match Level (M): Consider all **single** matches in the last month
  - Reduced Match Level ( $M^-$ ): consider at most 30 **single** matches in the last month
- Why need M &  $M^-$  ?
  - There are players with 200 matches a month, and others with 5 matches a month
  - If we don't tackle this imbalance, classifier would learn only from players with many matches, “hiding” signal of players with less matches
  - By putting a limit (e.g., 30 matches) we can obtain higher generalization

# Results – Correlation (overview)

Table 8: Significant Correlations at different  $p$ -values in our three datasets. Each column reports a personal attribute in  $\mathcal{A}$ . Rows denote how many features in each dataset (either  $\mathcal{M}$ ,  $\overline{\mathcal{M}}$  or  $\mathcal{P}$ ) achieve  $p$  below the target  $\alpha$  (i.e., the correlations are statistically significant).

<i>Dataset</i>	<i>Metric</i>	$\alpha$	gend.	age	occ.	purch.	extr.	agree.	consc.	neur.	open.
$\mathcal{M}$	Cram.	<0.01	17	17	15	18	13	18	17	16	13
	Cram.	0.05	18	19	15	18	14	19	18	19	14
	Cram.	0.1	18	19	17	19	15	19	19	19	16
	Spear.	0.01	–	88	–	51	44	52	22	70	36
	Spear.	0.05	–	95	–	65	57	59	35	85	50
	Spear.	0.1	–	99	–	73	62	67	43	87	59
$\overline{\mathcal{M}}$	Cram.	<0.01	16	12	12	11	15	10	10	14	8
	Cram.	0.05	18	17	18	15	17	11	14	15	11
	Cram.	0.1	18	17	18	15	18	14	15	20	13
	Spear.	0.01	–	95	–	43	53	38	25	60	27
	Spear.	0.05	–	104	–	63	65	54	40	82	47
	Spear.	0.1	–	108	–	69	73	64	53	90	58
$\mathcal{P}$	Cram.	<0.01	2	1	2	1	0	0	0	1	0
	Cram.	0.05	3	3	3	1	0	0	1	1	0
	Cram.	0.1	4	3	3	1	0	0	1	2	1
	Spear.	0.01	–	69	–	11	13	2	0	2	0
	Spear.	0.05	–	97	–	16	27	13	8	22	4
	Spear.	0.1	–	110	–	26	47	26	16	44	14

# Results – Correlation (detail)



# Results – Impact: Simple AIA (Aggregated data)

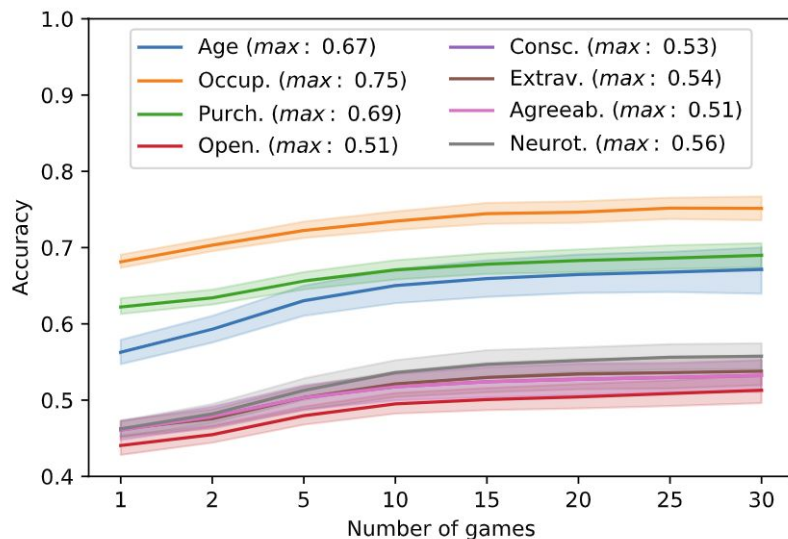
**Table 3: Impact of the *simple* AIA (based on  $\mathcal{P}$ ) as measured by the F1-score. Rows report the attributes and columns our ML models (boldface denotes the best model for a given attribute).**

	<i>LR</i>	<i>DT</i>	<i>RF</i>	<i>NN</i>	<i>Dummy</i>
gender	64.97 $\pm$ 10.9	59.71 $\pm$ 12.7	50.91 $\pm$ 5.33	<b>67.24<math>\pm</math>13.4</b>	51.62 $\pm$ 10.9
age	40.47 $\pm$ 6.30	39.38 $\pm$ 8.76	<b>44.08<math>\pm</math>6.17</b>	28.06 $\pm$ 7.59	32.21 $\pm$ 5.70
occup.	53.23 $\pm$ 7.22	47.44 $\pm$ 8.34	56.08 $\pm$ 7.88	<b>59.89<math>\pm</math>7.15</b>	43.76 $\pm$ 9.56
purch.	32.05 $\pm$ 10.1	31.74 $\pm$ 4.53	<b>34.40<math>\pm</math>8.20</b>	32.17 $\pm$ 7.19	31.20 $\pm$ 6.26
open.	28.94 $\pm$ 5.94	<b>40.76<math>\pm</math>6.80</b>	32.6 $\pm$ 7.77	30.89 $\pm$ 7.60	29.59 $\pm$ 2.04
consc.	26.52 $\pm$ 5.65	33.87 $\pm$ 8.78	<b>34.27<math>\pm</math>5.60</b>	23.83 $\pm$ 8.18	33.23 $\pm$ 8.94
extrav.	30.15 $\pm$ 7.53	36.16 $\pm$ 5.14	<b>36.49<math>\pm</math>5.56</b>	28.59 $\pm$ 5.95	32.27 $\pm$ 7.01
agreeab.	29.46 $\pm$ 6.29	<b>34.11<math>\pm</math>8.58</b>	33.68 $\pm$ 6.25	24.54 $\pm$ 9.43	33.39 $\pm$ 7.35
neurot.	32.38 $\pm$ 6.56	<b>40.76<math>\pm</math>6.80</b>	32.6 $\pm$ 7.74	31.6 $\pm$ 8.30	30.07 $\pm$ 4.46

# Results – Impact: Sophisticated One-to-One AIA

**Idea:** Build a match-based classifier, and use more matches to predict user's info

**Method:** Majority voting considering multiple matches



**Fig. 5: Impact of Sophisticated AIA.** The inference is done after post-processing the predictions of the ML model over multiple matches of the same targeted player (x-axis). The accuracy (y-axis) for all attributes (lines) increases as more matches are considered.



# Results – Impact: Indiscriminate Many-to-Many AIA

**Idea:** The attacker is fine with “not completely wrong” predictions

**Method:** Consider both first and second predictions as correct

**Table 6: Indiscriminate ‘many-to-many’ AIA (mid column). Compared to the baseline (cf. Fig. 5), the accuracy substantially increases.**

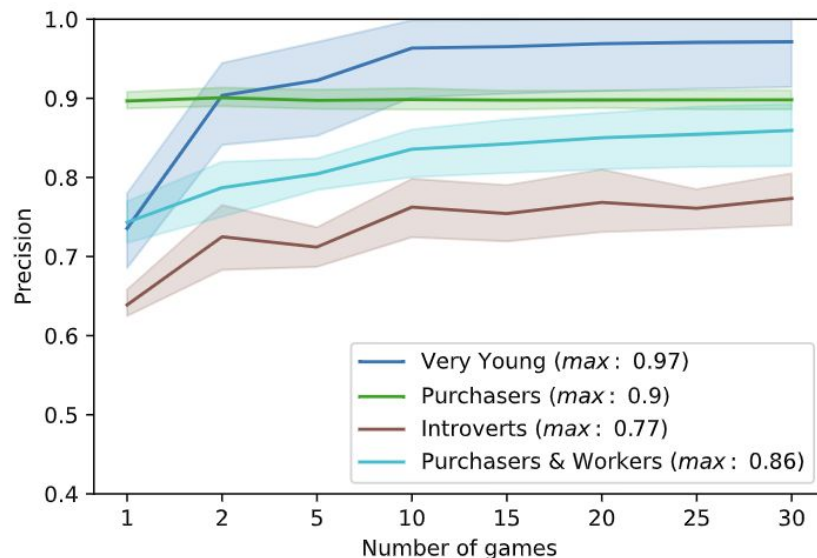
	Sophisticated AIA (30 matches)	Indiscriminate AIA (30 matches)	Improvement
age	67.15 $\pm$ 6.87	<b>89.15</b> $\pm$ 4.66	+22.00%
purch.	68.99 $\pm$ 3.81	<b>96.13</b> $\pm$ 2.86	+27.14%
open.	51.30 $\pm$ 3.87	<b>77.86</b> $\pm$ 3.39	+26.56%
consc.	53.24 $\pm$ 4.88	<b>80.19</b> $\pm$ 4.12	+26.95%
extrav.	53.78 $\pm$ 3.90	<b>81.51</b> $\pm$ 4.40	+27.73%
agreeab.	50.71 $\pm$ 4.65	<b>76.84</b> $\pm$ 5.59	+26.13%
neurot.	55.74 $\pm$ 3.88	<b>80.64</b> $\pm$ 4.02	+24.90%



# Results – Impact: Targeted Many-to-One AIA

**Idea:** The attacker wants to be precise in finding a target, not in finding all of them

**Method:** Train and validate models to reach high precision



**Fig. 6:** Targeted ‘many-to-one’ AIA. We train our ML models by maximizing the *precision* on a single targeted class. Such AIA are very effective after analyzing 10 matches for each player in the test-set.

# So What Now?

- **Hard counters? Nope!**
  - The entire E-sport ecosystem would be disrupted
- **Compromise? Yes!**
  - The users should be informed that having their in-game statistics to be publicly accessible by TW exposes them to AIA
  - Access control rules
  - Turn TW into social networks
  - All of these require effort and collaboration between VG and TW (not easy!)

# So What Now?

- **What about other games?** Many E-sports share the same ecosystem with Dota2
  - AIA are theoretically possible also in other VG, but a correlation has to be found first

Table 7: Overview of E-Sports VG. Numbers are taken from various sources [17, 20, 32, 52, 59].

	Release Year	Genre	Monthly Players	Concurrent Players Avg	Playtime Avg	Age Range (PEGI rec.)	Tournament Revenue	Exemplary TW	Replay System	Max Players per Lobby
<i>League of Legends</i>	2009	MOBA	127 M	700 K	832 H	11–50 (12+)	\$93 M	lolprofile.net	Yes	10
<i>CS:GO</i>	2012	FPS	34 M	560 K	611H	13–40 (18+)	\$134 M	csgostats.gg	Yes	18
<i>Rocket League</i>	2016	Sport	90 M	25 K	315 H	6–35 (3+)	\$18 M	rltracker.pro	Yes	8
<i>Fortnite</i>	2017	Battle Royale	270 M	4 M	1800 H	6–54 (12+)	\$121 M	fortnitetracker.com	Yes	100
<i>PUBG</i>	2018	Battle Royale	510 M	200 K	356 H	12–55 (16+)	\$45 M	pubg.op.gg	Yes	100
<i>Apex Legends</i>	2019	Battle Royale	118 M	195 K	91 H	8–37 (16+)	\$10 M	apex.tracker.gg	No	60
DOTA2	2013	MOBA	3.7 M	450 K	1700H	12–50 (12+)	\$283 M	opendota.com	Yes	10

- **We sent an email to Valve** to inform them of such vulnerability.
  - We are unsure about whether they will take any action in the short-term

Thank you!

Questions?

**Pier Paolo Tricomi**

*tricomi@math.unipd.it*