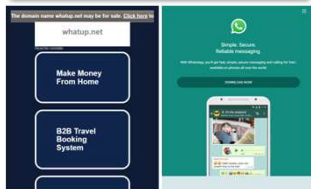


Research context and motivation

- One common attack is **domain squatting**, which occurs when attackers register **perceptively confusing domain names** aiming at tricking visitors into them.
- Sound-squatting** has gaining traction with the advent of **smart speakers and voice-assistants**.
- The **state-of-art** in detection uses **statically built lists of homophones**.
- We **hypothesize** Artificial Intelligence can produce **more comprehensive** sound-squatting candidates and be used as an automatic method for **sound-squatting generation and detection**.

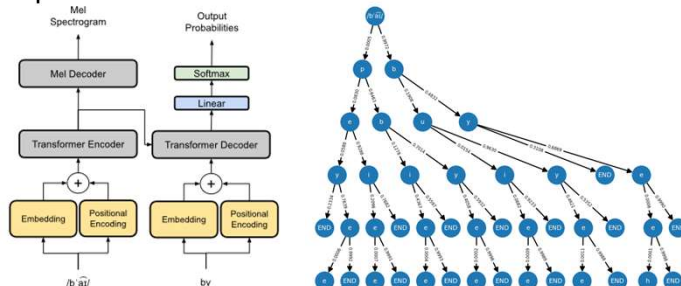
Domain Squatting Types

bit qoutube.com
typo yotube.com
sound utube.com
combo videoyoutube.com
homograph yovtube.com



Adopted methodologies

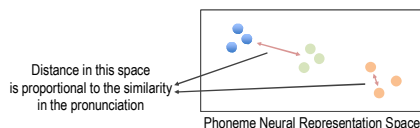
- We employ **Transformer models** to translate from **International Phoneme Alphabet (IPA) to English (from Phoneme to Grapheme)**.
- The model learns how to **map the IPA tokens into an audio dependable representation**.



- The methodology for **searching domain squatting** consists of **generating candidates** from **high profile domains** and **actively verifying** if these domains are being abused.

Addressed research questions/problems

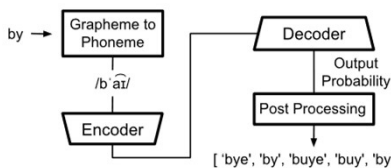
- How to **proactively generate** sound-squatting candidates **with AI**?
- Can we learn a **phoneme representation** that considers **sound features**?



- Can we automatically **validate** the homophone **generation**?
- Can we generate **more than one candidate** from the same target **maintaining quality**?

Novel contributions

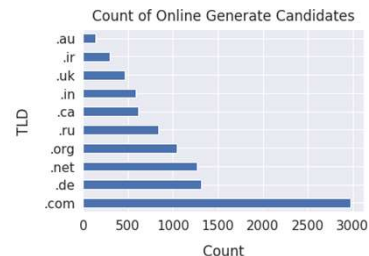
- We propose a **novel method** for **generation of sound-squatting candidates**.



- We find a **neural representation** to mapping from International Phoneme Alphabet (IPA), English written symbols (**grapheme**) and **actual pronunciations**.
- We propose a **Post Processing** method that extracts from the Transformers not only the **best** output, but **also acceptable** outputs accessing the quality and introducing variability.
- We propose a **proactive search using AI-generated candidates** extending the start-of-art that uses statically built lists retrieved from known homophones.

Preliminary Results

Domain	Candidates
google	googal, gougall, gougale, googall, gugal
youtube	utube, uteube, yutube, yootube
facebook	phasebook, phacebook, facebooke



- To validate the method, from a list of **2279 set of homophones** we **generated 66017 candidates** out of which **84.17% (4885)** are in the set of **known 5804 homophones**.

Future work

- We will **extend** our methodology for **other languages** which is not trivial due to **phoneme gaps**.
- We also plan to **inspect the occurrence** of sound-squatting to **other contexts**, such as: python packages and smart speaker voice commands.



List of attended classes

- 01UMNRV - Advanced deep Learning (didattica di eccellenza) (15/06/2021, 40.00)
- 01UJBRV - Adversarial training of neural networks (03/06/2021, 25.00)
- 01TRARV - Big data processing and programming (08/03/2021, 33.33)
- 02LWHRV - Communication (06/08/2022, 6.67)
- 01UJTUV - Control and data acquisition automation in scientific experiments (15/02/2021, 16.67)
- 01QTEIU - Data mining concepts and algorithms (01/02/2021, 33.33)
- 01SHMRV - Entrepreneurial Finance (07/08/2022, 6.67)
- 01UNVRV - Navigating the hiring process: CV, tests, interview (05/08/2022, 2.67)
- 01UNYRV - Personal branding (05/08/2022, 1.33)
- 01RISRV - Public speaking (05/09/2022, 6.67)
- 01SYBRV - Research integrity (05/08/2022, 6.67)
- 01SWQRV - Responsible research and innovation, the impact on social challenges (06/08/2022, 6.67)
- 01TSBRV - Scienza dei dati applicata alle reti complesse (08/09/2022, 26.67)
- 02RHORV - The new Internet Society: entering the black-box of digital innovations (06/08/2022, 8.00)
- 01UNXRV - Thinking out of the box (05/08/2022, 1.33)
- 01SWPRV - Time management (05/09/2022, 2.67)

Soft Skills Hours: 37/40
Hard Skills Score: 175/200

Submitted and published works

- Valentim, R. V., Comarela, G., Park, S. and Sáez-Trumper, D. (2021). Tracking Knowledge Propagation Across Wikipedia Languages. Proceedings of the International AAAI Conference on Web and Social Media, 15(1), 1046-1052. <https://ojs.aaai.org/index.php/ICWSM/article/view/18128>
- Valentim, R., Drago, I., Trevisan, M., Cerutti, F., and Mellia, M. (2021). Augmenting Phishing Squatting Detection with GANs. Proceedings of the CoNEXT Student Workshop, 3-4. <https://doi.org/10.1145/3488658.3493787>
- Valentim, R., Drago, I., Cerutti, F., and Mellia, M. (2022). AI-based Sound-Squatting Attack Made Possible. 2022 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW), 448-453. <https://doi.org/10.1109/EuroSPW55150.2022.00053>
- Valentim, R., Drago, I., Trevisan, M., and Mellia, M. (2022). URLGEN - Towards Automatic URL Generation Using GANs. IEEE Transactions on Network and Service Management (TNSM). (Submitted)