

A Testbed for Voice Assistant Traffic Fingerprinting

Master Thesis Presentation

Milan van Zanten

University of Basel

21.03.2024

Outline

1. Voice Assistants
2. Traffic Fingerprinting
3. Testbed
4. Results
5. Demo

Ask questions any time!

A Testbed for *Voice Assistant* Traffic Fingerprinting

Specifically, *Smart Speakers*



Echo Dot

Amazon Alexa



HomePod Mini

Siri



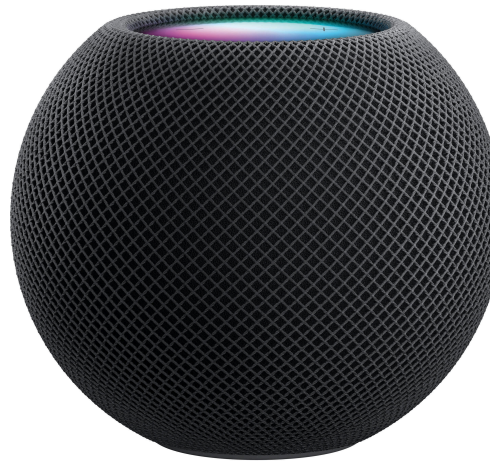
Google Home Mini

Google Assistant

Specifically, *Smart Speakers*



Echo Dot
Amazon Alexa



HomePod Mini
Siri



Google Home Mini
Google Assistant

none (J °□°) J ˘ 11

There are concerns...

- Usually located where sensitive conversations take place
- Necessarily always listening
 - Misactivations
- Used to control smart home devices (e.g. door locks)
- No authentication*

* Voice recognition is still insecure.

There are concerns...

- Usually located where sensitive conversations take place
- Necessarily always listening
 - Misactivations
- Used to control smart home devices (e.g. door locks)
- No authentication*

About 40% of households in the U.S. own a smart speaker.

* Voice recognition is still insecure.

Active:

- Malicious activations
- Similar pronunciations, “skill squatting”
 - (e.g. “Boil an egg” → “Boyle an egg”)¹

Passive:

- Traffic Fingerprinting

¹D. Kumar et al., “*Skill Squatting Attacks on Amazon Alexa*”, August 2018, Available: <https://www.usenix.org/conference/usenixsecurity18/presentation/kumar>

Active:

- Malicious activations
- Similar pronunciations, “skill squatting”
 - (e.g. “Boil an egg” → “Boyle an egg”)²

Passive:

- **Traffic Fingerprinting**

²D. Kumar et al., “Skill Squatting Attacks on Amazon Alexa”, August 2018, Available: <https://www.usenix.org/conference/usenixsecurity18/presentation/kumar>

A Testbed for Voice Assistant Traffic Fingerprinting

“[SSL] traffic analysis aims to recover confidential information about protection sessions by examining unencrypted packet fields and unprotected packet attributes. For example [...] the volume of network traffic flow”

— Wagner and Schneier³

³D. Wagner and B. Schneier, “*Analysis of the SSL 3.0 Protocol*”, November 1996, Available: <https://dl.acm.org/doi/10.5555/1267167.1267171>

“[SSL] traffic analysis aims to recover confidential information about protection sessions by examining unencrypted packet fields and **unprotected packet attributes**. For example [...] the volume of network traffic flow”

— Wagner and Schneier⁴

⁴D. Wagner and B. Schneier, “*Analysis of the SSL 3.0 Protocol*”, November 1996, Available: <https://dl.acm.org/doi/10.5555/1267167.1267171>

“[SSL] traffic analysis aims to recover confidential information about protection sessions by examining unencrypted packet fields and **unprotected packet attributes**. For example [...] the volume of network traffic flow”

— Wagner and Schneier⁵

... packet direction, timing, and more

⁵D. Wagner and B. Schneier, “*Analysis of the SSL 3.0 Protocol*”, November 1996, Available: <https://dl.acm.org/doi/10.5555/1267167.1267171>

1996 Wagner and Schneier¹, coined SSL traffic analysis

1998 Cheng and Avnur, website traffic analysis

website fingerprinting (WF)...

2016 Abe and Goto, deep learning WF

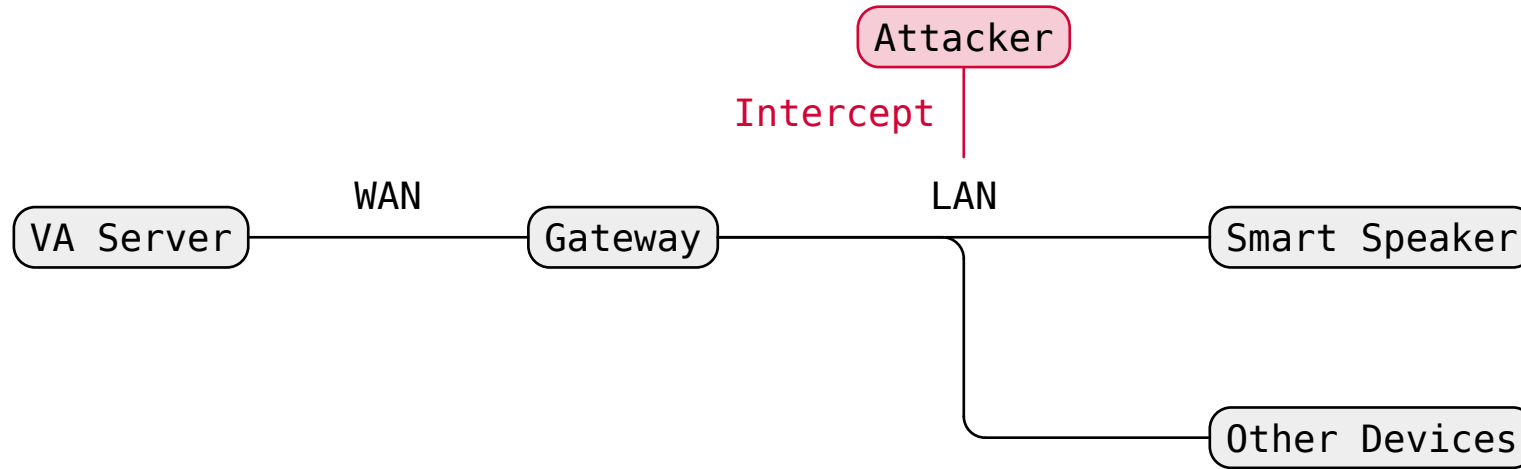
2019 Kennedy et al., apply WF techniques to voice assistants (VA)

2020 Wang et al., deep learning VA fingerprinting

2022 Mao et al., temporal features

2023 Ahmed, Sabir and Das, invocation detection

¹Timeline references can be found at the end of the presentation.



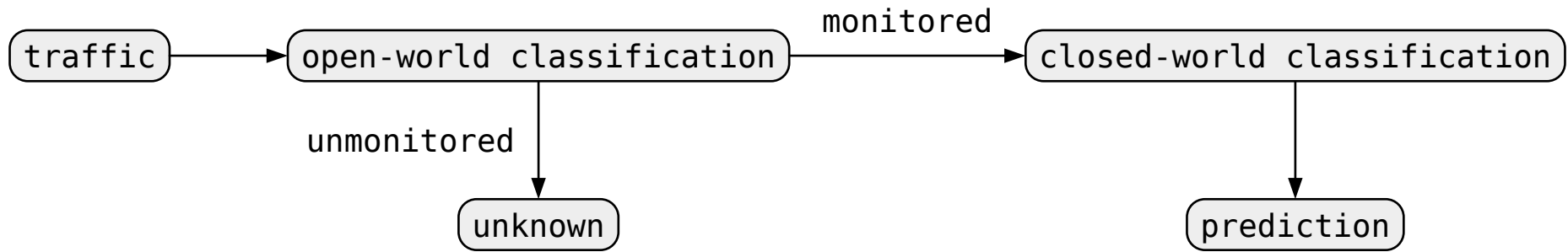
1. The attacker can intercept traffic from smart speaker
2. The attacker knows the smart speaker address
3. The attacker knows the type of smart speaker used
4. The attacker knows the beginning and end of an interaction

- Fixed list of monitored voice commands
- Traffic is considered to come from one of the monitored commands
- Multiclass classification

Predicts which command was used.

- Traffic can also come from new, unmonitored commands
- Binary-classification

Predicts whether traffic is from monitored or unmonitored command.



A **Testbed** for Voice Assistant Traffic Fingerprinting

Website Fingerprinting:

- Requires a large amount of data
- Data collection usually via program making requests
- Only dependent on network environment
- Fast

Voice Command Fingerprinting:

Website Fingerprinting:

- Requires a large amount of data
- Data collection usually via program making requests
- Only dependent on network environment
- Fast

Voice Command Fingerprinting:

- Requires a large amount of data

Website Fingerprinting:

- Requires a large amount of data
- Data collection usually via program making requests
- Only dependent on network environment
- Fast

Voice Command Fingerprinting:

- Requires a large amount of data
- Interaction by speaking

Website Fingerprinting:

- Requires a large amount of data
- Data collection usually via program making requests
- Only dependent on network environment
- Fast

Voice Command Fingerprinting:

- Requires a large amount of data
- Interaction by speaking
- Hampered by environment noise

Website Fingerprinting:

- Requires a large amount of data
- Data collection usually via program making requests
- Only dependent on network environment
- Fast

Voice Command Fingerprinting:

- Requires a large amount of data
- Interaction by speaking
- Hampered by environment noise
- Slow and inefficient

Website Fingerprinting:

- Requires a large amount of data
- Data collection usually via program making requests
- Only dependent on network environment
- Fast

Voice Command Fingerprinting:

- Requires a large amount of data
- Interaction by speaking
- Hampered by environment noise
- Slow and inefficient

→ Sophisticated testbed

“The content of voice commands may vary from date to date;
therefore, more efficient data collection tools need to be developed.”
— Mao et al.¹

¹Jianghan Mao et al., “A novel model for voice command fingerprinting using deep learning”, March 2022, Available: <https://doi.org/10.1016/j.jisa.2021.103085>

- Sound isolation
 - Isolated box
 - Separate speaker/microphone
- Efficiency
 - Every second saved per interaction means hours saved when collecting dozens of thousands interactions
 - Dynamic interaction length by listening for silence
- Robustness
 - Autonomously reset VA if error occurs
 - Monitoring system

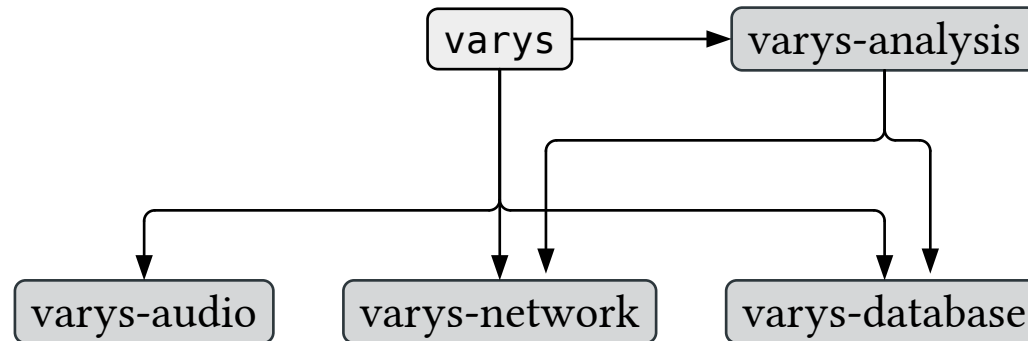
varys The main executable combining all modules into the final system.

varys-analysis Analysis of data collected by varys.

varys-audio Recording audio and the TTS and STT systems.

varys-database Abstraction of the database system where interactions are stored.

varys-network Collection of network traffic, writing and parsing of .pcap files.



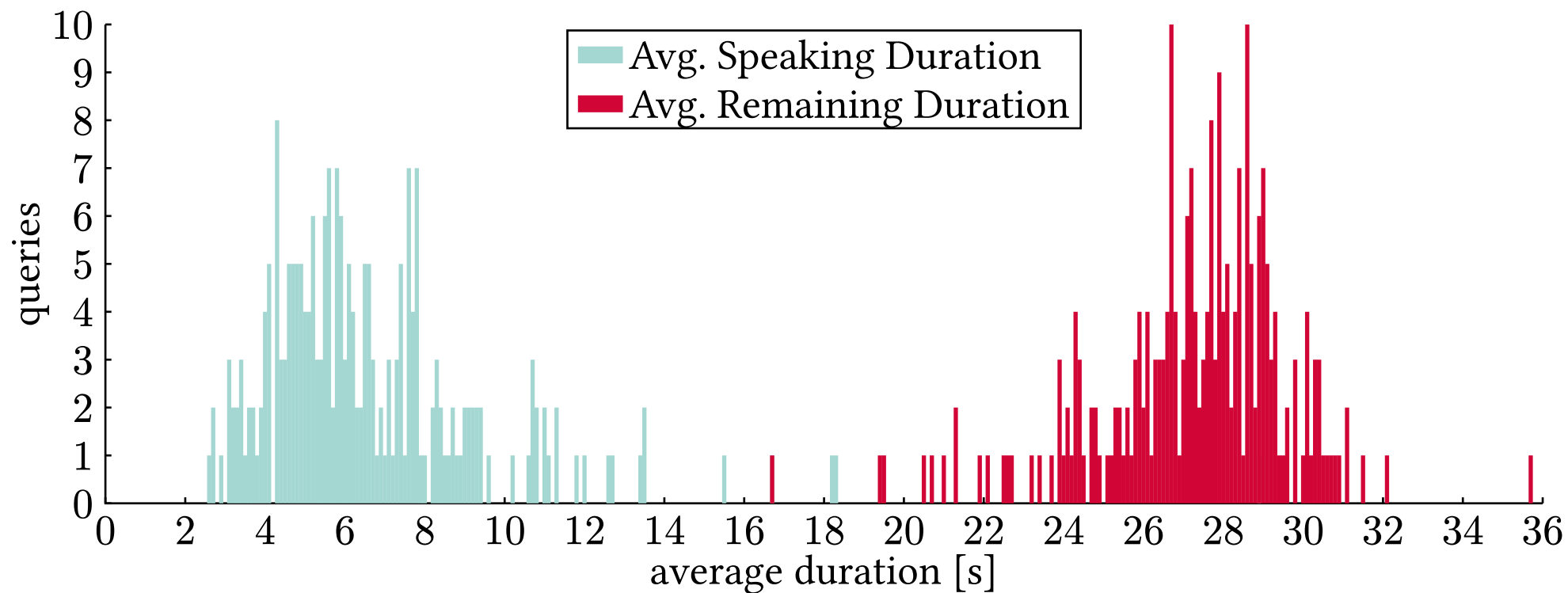
Results

~800h, ~70'000 interactions

large 227 queries, 140 interactions each

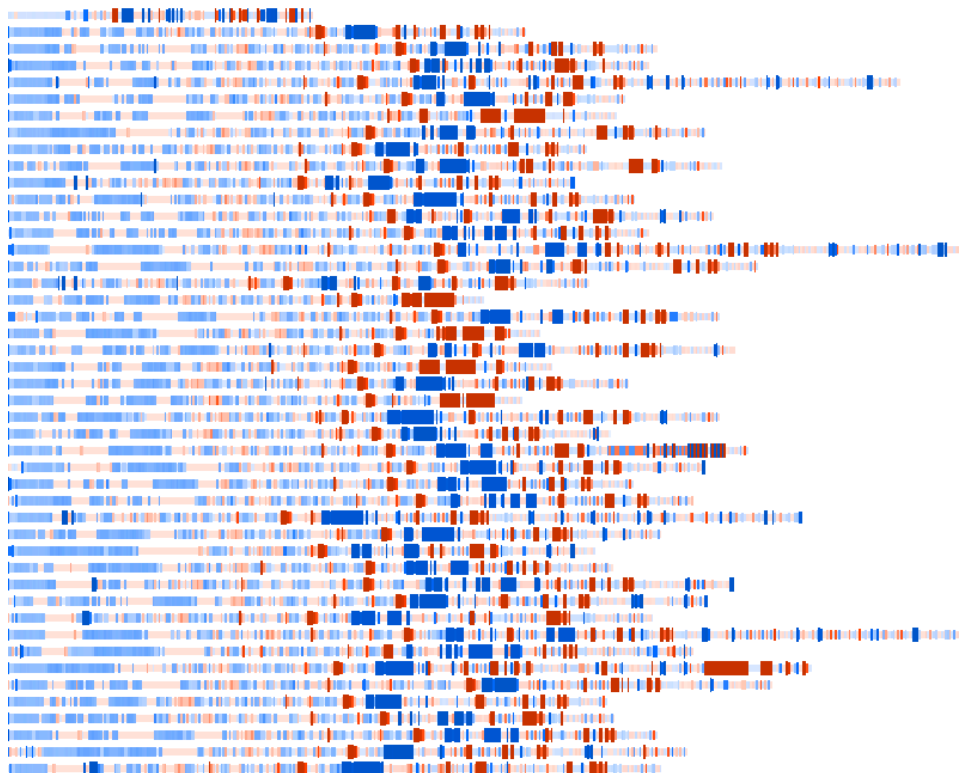
small 13 queries, 2400 interactions each

binary “*Call John Doe*” and “*Call Mary Poppins*”, 1500 interactions each



Traffic Trace Examples

Results



“Any missed calls?”

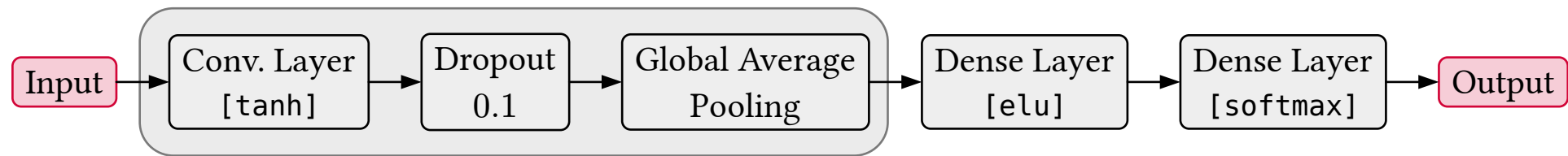


“What day was 90 days ago?”

Feature extraction (packet size $s \in [0, 1500]$ and direction $d \in \{0, 1\}$):

$$(s, d) \rightarrow (-1)^d \cdot \frac{s}{1500}$$

CNN adapted from Wang et al.¹:



¹Chenggang Wang et al., “Fingerprinting Encrypted Voice Traffic on Smart Speakers with Deep Learning”, May 2020, Available: <https://doi.org/10.1145/3395351.3399357>

Accuracy on test sets:

large ~40.40% (random choice ~0.44%)

small ~86.19% (random choice ~7.69%)

binary ~71.19% (random choice 50%)

Demo

```
./target/release/varys -i ap1 analyse demo data/ml/test_5_13\ queries_0.86 f4:34:f0:89:2d:75
```

“Hey Siri, any missed calls?”

“Hey Siri, remind me to wash the car.”

It is unlikely this will work...

Timeline References

- D. Wagner and B. Schneier, “*Analysis of the SSL 3.0 Protocol*”, November 1996, Available: <https://dl.acm.org/doi/10.5555/1267167.1267171>
- H. Cheng and R. Avnur, “*Traffic Analysis of SSL Encrypted Web Browsing*”, 1998
- K. Abe and S. Goto, “*Fingerprinting Attack on Tor Anonymity using Deep Learning*”, August 2016, Available: <https://core.ac.uk/display/229876143>
- S. Kennedy et al., “*I Can Hear Your Alexa: Voice Command Fingerprinting on Smart Home Speakers*”, June 2019, Available: <https://doi.org/10.1109/CNS.2019.8802686>
- Chenggang Wang et al., “*Fingerprinting Encrypted Voice Traffic on Smart Speakers with Deep Learning*”, May 2020, Available: <https://doi.org/10.1145/3395351.3399357>
- Jianghan Mao et al., “*A novel model for voice command fingerprinting using deep learning*”, March 2022, Available: <https://doi.org/10.1016/j.jisa.2021.103085>
- D. Ahmed, A. Sabir, and A. Das, “*Spying through Your Voice Assistants: Realistic Voice Command Fingerprinting*”, August 2023, Available: <https://www.usenix.org/conference/usenixsecurity23/presentation/ahmed-dilawer>