

(Un-)Natural language processing

Defensive AI in practice

Mages

Monkeys

Bananas

Leopards

CrowdSec TLDR

- French Startup founded in 2020
- Team of ~35 people distributed over Europe
- Been there for a bit more than 2 years
- We produce an open source IDPS/WAF

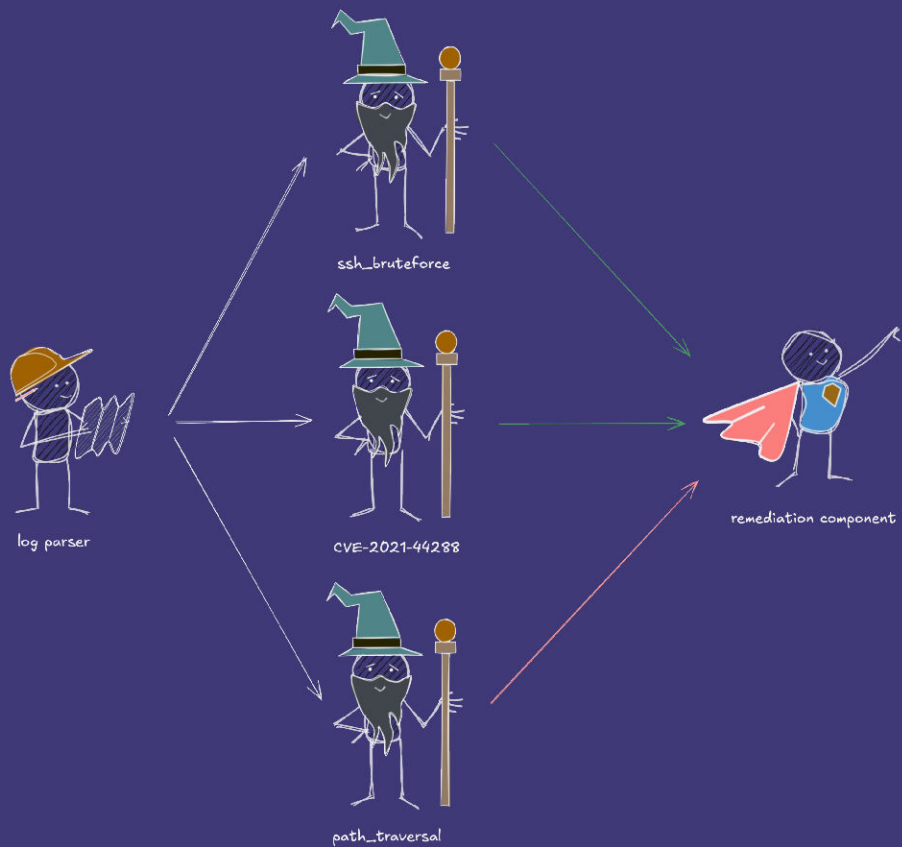


Security software TLDR

Lightweight agent that reads your logs/traffic

Logs/traffic are matched with rulesets/signatures

Detected intruders are sent to remediation



but...

Who teaches the mages to read?

The Update Problem

For a security tool to be effective, the detection signatures need to be updated continually and in a timely manner.

This incurs a cost both on our side and on the customer side.

On the other hand:

`/blog/the-new-air-jordan-is-here`

`/store/home`

`/store/search?query=air%20jorden`

`/store/search?query=air%20jordan`

`/store/search?_jinjadbg=.forName("java.lang.Runtime`

`/cart?product_id=374987&quantity=12`

`/checkout`

On the other hand:

/blog/the-new-air-jordan-is-here

/store/home

/store/search?query=air%20jorden

/store/search?query=air%20jordan

/store/search?_jinjadbg=.forName("java.lang.Runtime

/cart?product_id=374987&quantity=12

/checkout



Conclusion

If you had a monkey that would read all your logs, it could detect attacks by simply reporting anything that looks abnormal!

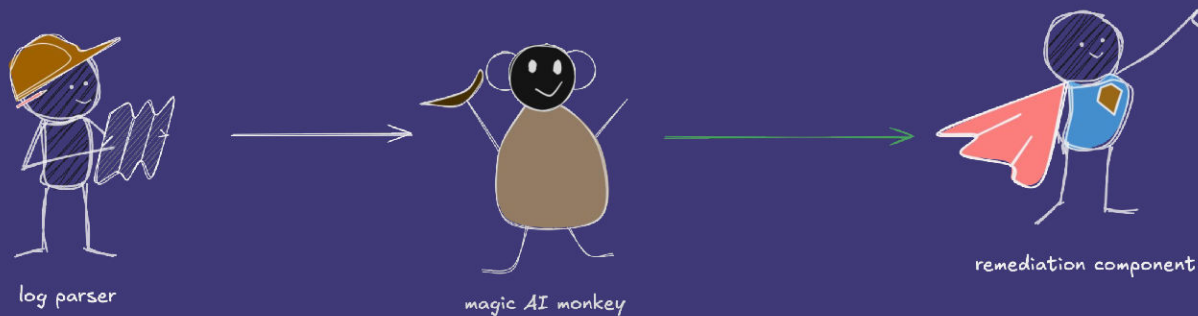
Conclusion

If you had a monkey that would read all your logs, it could detect attacks by simply reporting anything that looks abnormal!

=> This is a good usecase for AI

The Plan

Create a generic outlier detection algorithm to act as a first line of defense against novel threats.



Questions going into the project

- Does this famed AI monkey even exist?
- Can we transfer monkeys or do we have to train one per customer?
- Are these monkeys cost effective?

The trick

URIs have a known structure which we can use to turn a URI into something that looks like a sentence.

Afterwards we simply copy-paste the tricks we learned from natural language processing.

Example

“/store/search?query=air%20jordan”

Example

“/store/search?query=air%20jordan”

=> split on /, ?, = (url structure)

“store search query air%20jordan”

Example

“/store/search?query=air%20jordan”

=> split on /, ?, = (url structure)

“store search query air%20jordan”

=> split on “%20” (url encoding)

“store search query air jordan”

The method

Based on “sec2vec” paper by Gniewkowski et al.

Basically a two step process:

- Turn URI into a vector (embedding)
- Run a second ML model to classify embeddings into good and bad

Sec2Vec pipeline

URI → English sentence (base transform)

English → sequence of numbers (tokenizer)

seq. of numbers → vector (embedding)

vector → label (classifier)

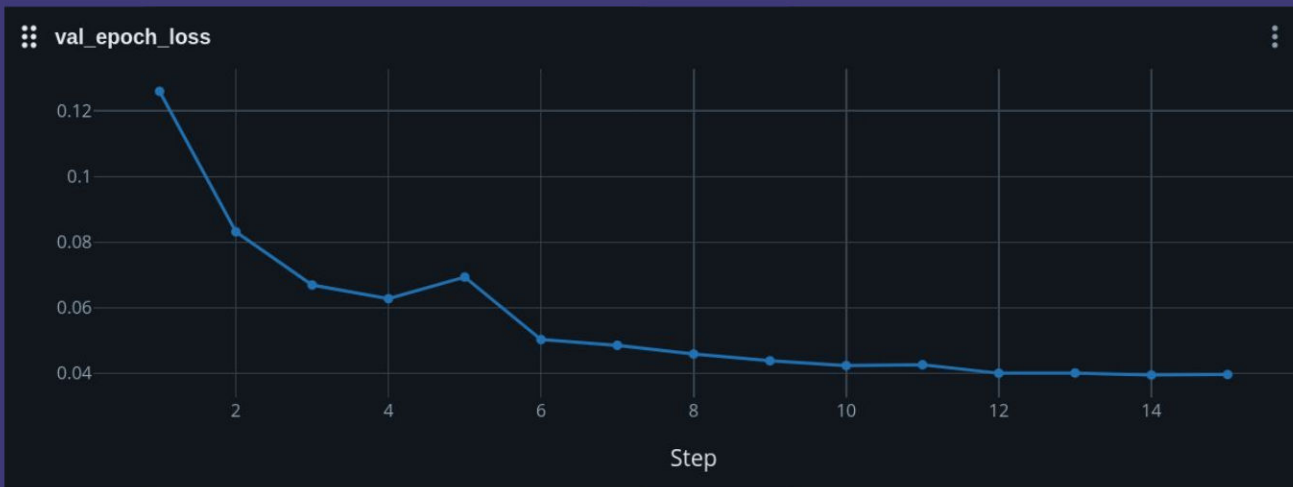
Results

Sec2Vec actually works!*



Results

Sec2Vec actually works!*



* on research data

The data problem

Due to various constraints (legal, privacy etc.) the logs publicly available to researchers are terrible

The data problem

Due to various constraints (legal, privacy etc.) the logs publicly available to researchers are terrible

In fact, this year there was a paper on how bad they generalize to real logs

Conferences > 2024 IEEE 9th European Sympos... ?

Bad Design Smells in Benchmark NIDS Datasets

Publisher: IEEE

Cite This



Robert Flood ; Gints Engelen ; David Aspinall ; Lieven Desmet [All Authors](#)

Trading NDAs for bananas

To fix the issue with the synthetic research dataset we reached out to our users and customers to send us their logs



Benefits on better data

Models trained on our dataset performed significantly better, dropping about 9/10 attacks used by a common vulnerability scanner

The new models also generalized better to new url structures, i.e. one model could serve multiple users

Conclusions

- **Researchers have pretty good ideas**
- **Researchers have pretty bad data**
- **With better data, language models have serious potential in network defense**

Bonus Chapter: Performance

To illustrate the trade-offs faced by security teams, consider a standard website service with some form of security software attached

We can differentiate between two modes of security operations: in-band and out-of-band

In-band system



User sees content only after they pass security

Out-of-band system



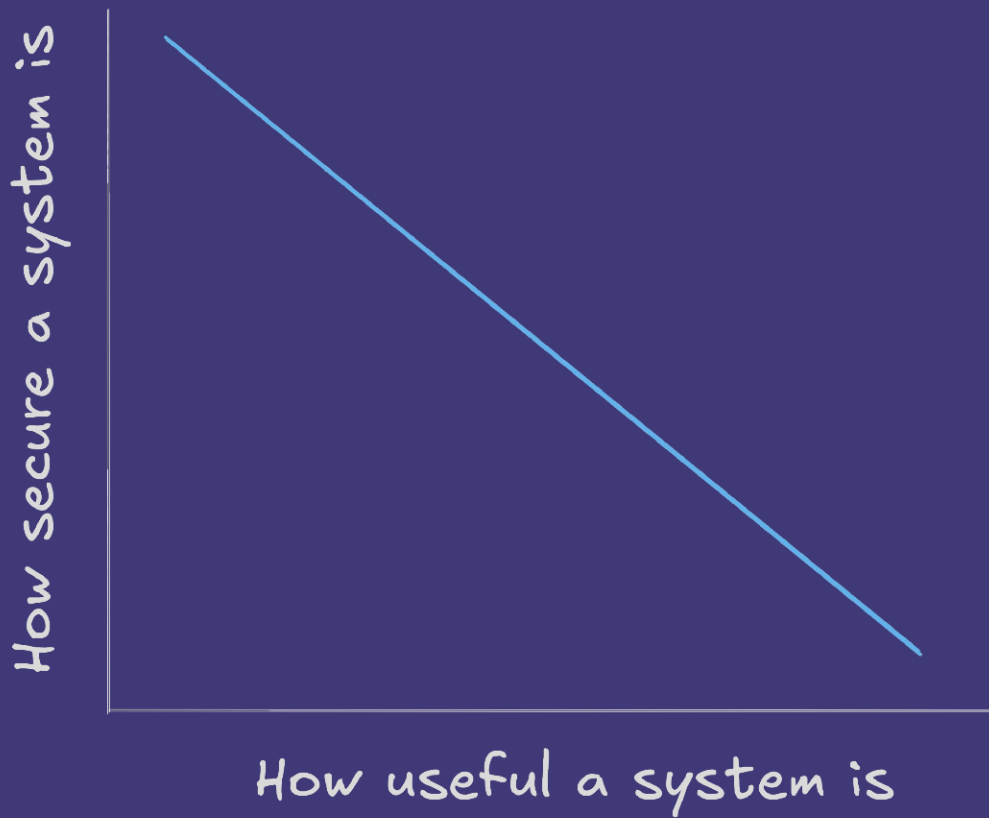
User can connect before the security system cleared them



If out-of-band is so insecure...

Why would anyone use an out-of-band security solution?

The security tradeoff



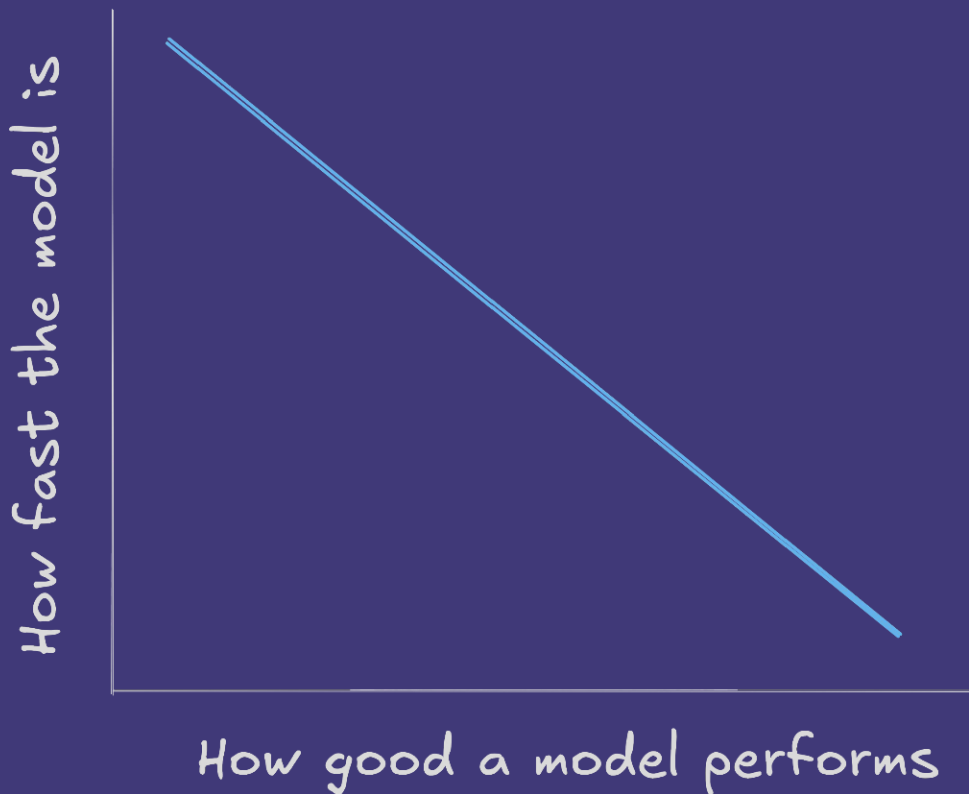
Performance considerations

In an in-band system, a user will see higher time-to-first-byte, an important business objective

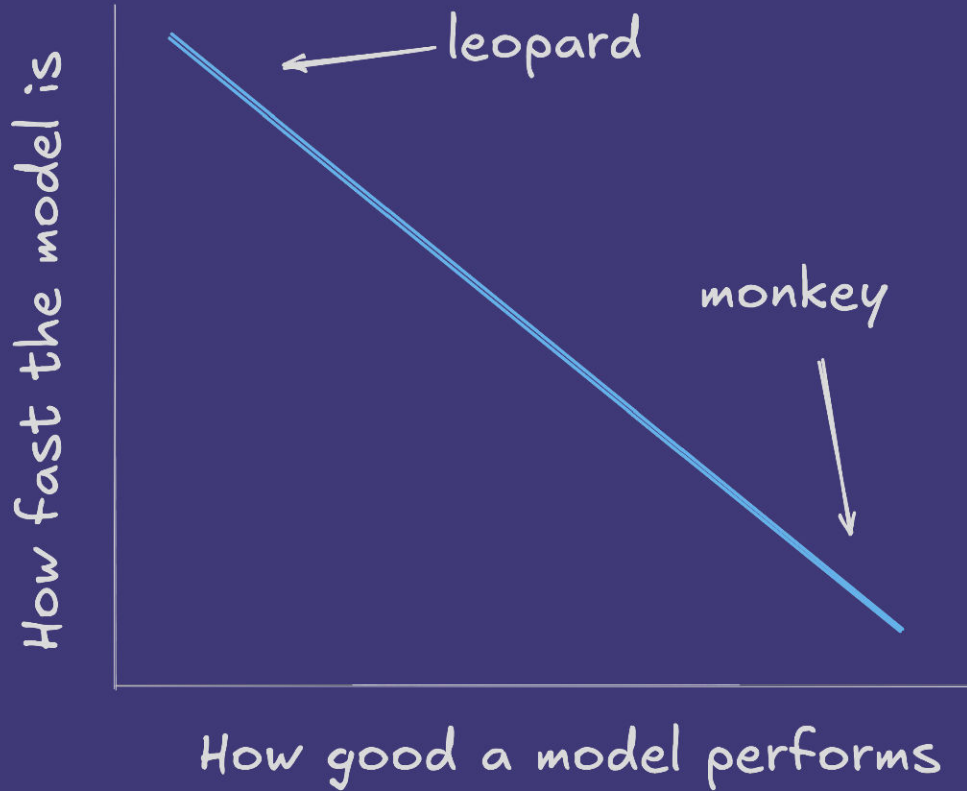
In an out-of-band system, an attacker might get in before the security system has time to respond

Such competing objectives we also have to worry about with the ML models

The Deep Learning tradeoff



The monkey leopard tradeoff



Performance considerations

Big ML models (monkeys) require specialized hardware to run fast which is unnecessary cost factor for a webserver

There are optimizations that can be done, such as quantization or decreasing model size (leopards) but these come at an effectiveness cost

Conclusions

- Researchers have pretty good ideas
- Researchers have pretty bad data
- With better data, language models have serious potential in network defense
- Feasibility of running big ML models still an open question

Answers coming out of the project

- Does this famed AI monkey even exist?
- Can we transfer monkeys or do we have to train one per customer?
- Are these monkeys cost effective?

Answers coming out of the project

- Does this famed AI monkey even exist?
=> Yes
- Can we transfer monkeys or do we have to train one per customer?
- Are these monkeys cost effective?

Answers coming out of the project

- Does this famed AI monkey even exist?

⇒ Yes

- Can we transfer monkeys or do we have to train one per customer?

⇒ Yes, if we feed them enough bananas

- Are these monkeys cost effective?

Answers coming out of the project

- Does this famed AI monkey even exist?

⇒ Yes

- Can we transfer monkeys or do we have to train one per customer?

⇒ Yes, if we feed them enough bananas

- Are these monkeys cost effective?

⇒ No, but we can train leopards instead



Next steps

We are open-sourcing the models and the plugin to run sec2vec within the crowdsec software soon-ish (Q4/Q1)

Repo:

github.com/crowdsecurity/crowdsec

Contact:

datascience@crowdsec.net

emanuel@crowdsec.net