# Yarlx: Scalable YARA-based Malware Intelligence

Michael Brengel, Christian Rossow

CISPA Helmholtz Center for Information Security
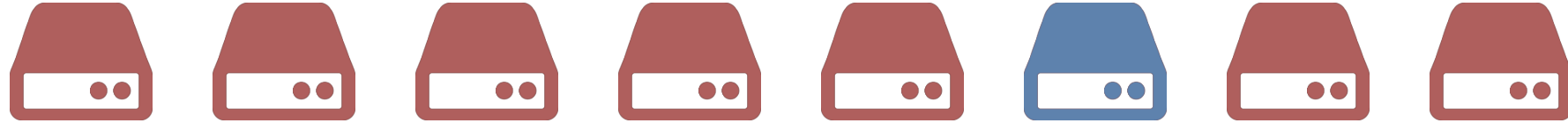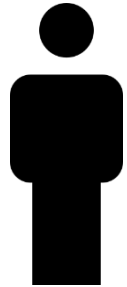
https://github.com/mbrengel/yarix
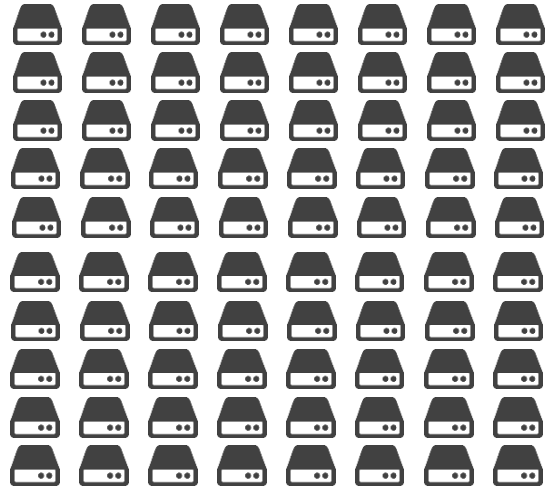
ARTIFACT
EVALUATED

usenix
ASSOCIATION

PASSED

Which malware samples contain the string "\ScreenBlaze.exe"?

YarIx: Scalable YARA-based Malware Intelligence
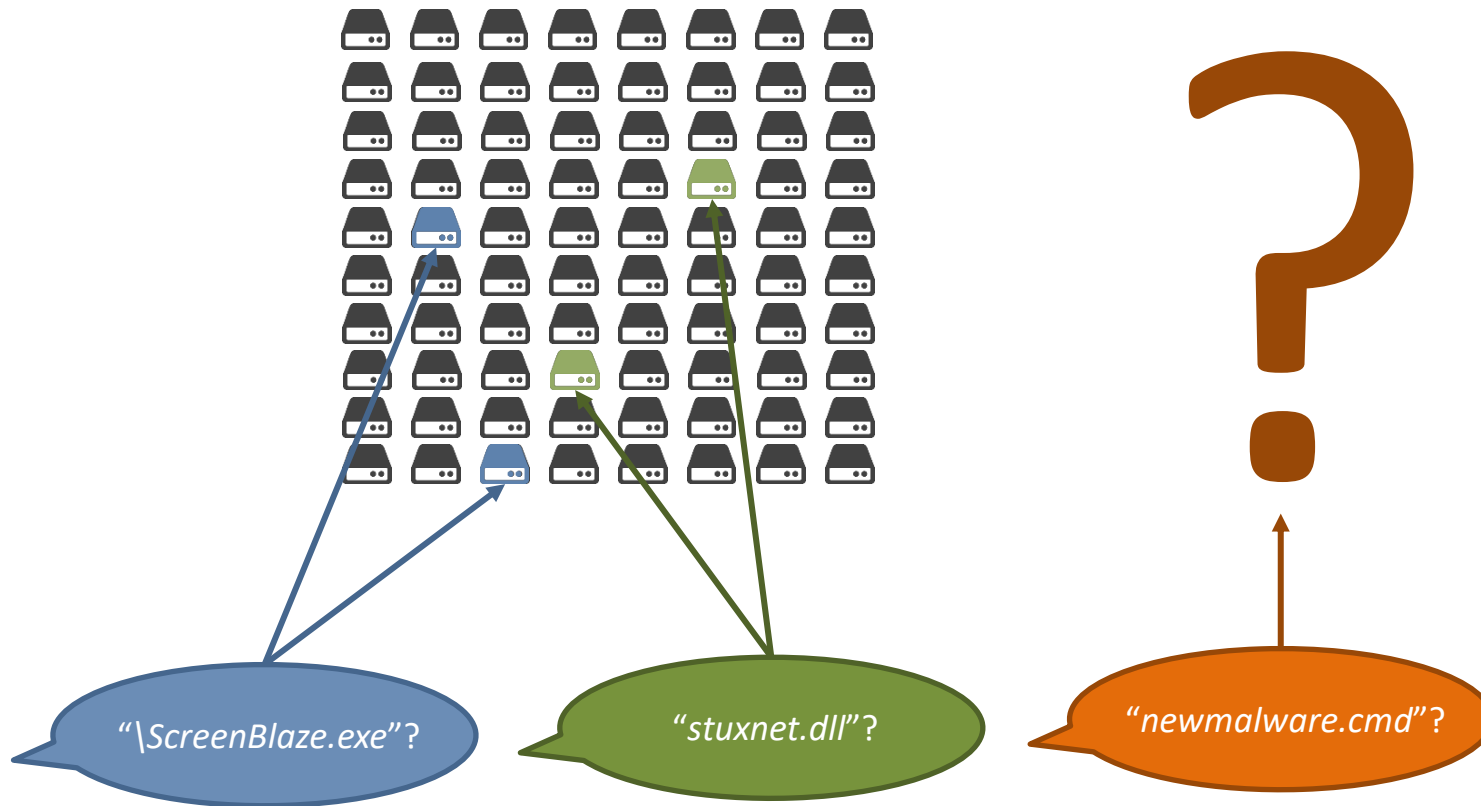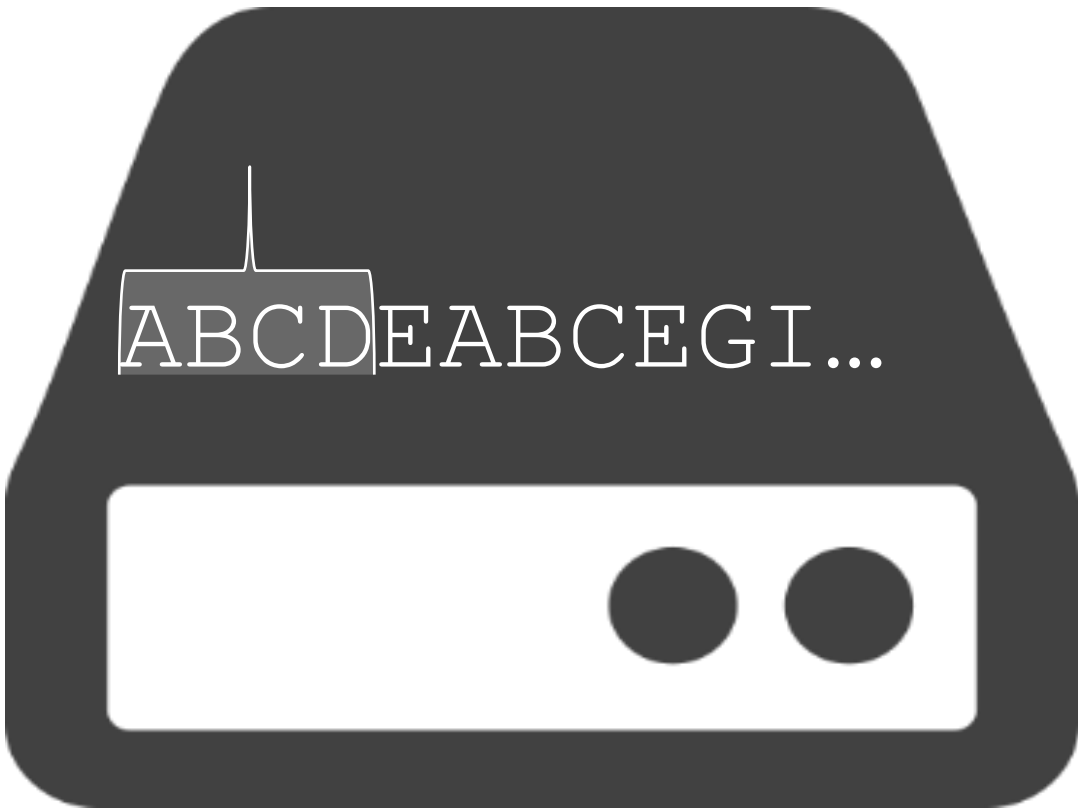
How can we use YARA *efficiently* in *large* malware databases?

Which malware samples contain the string "*\ScreenBlaze.exe*"?

YarIx: Scalable YARA-based Malware Intelligence

# YarIx: Reverse Malware Index

ABCDEABCEGI...

| 4-gram | Posting Lists |
|--------|---------------|
| ABCB | |
| ABCD | |
| ABCE | |
| ABCF | |
| BCDE | |
| CDEA | |
| CEGI | |
| DEAB | |

YarIx: Scalable YARA-based Malware Intelligence

CISPA
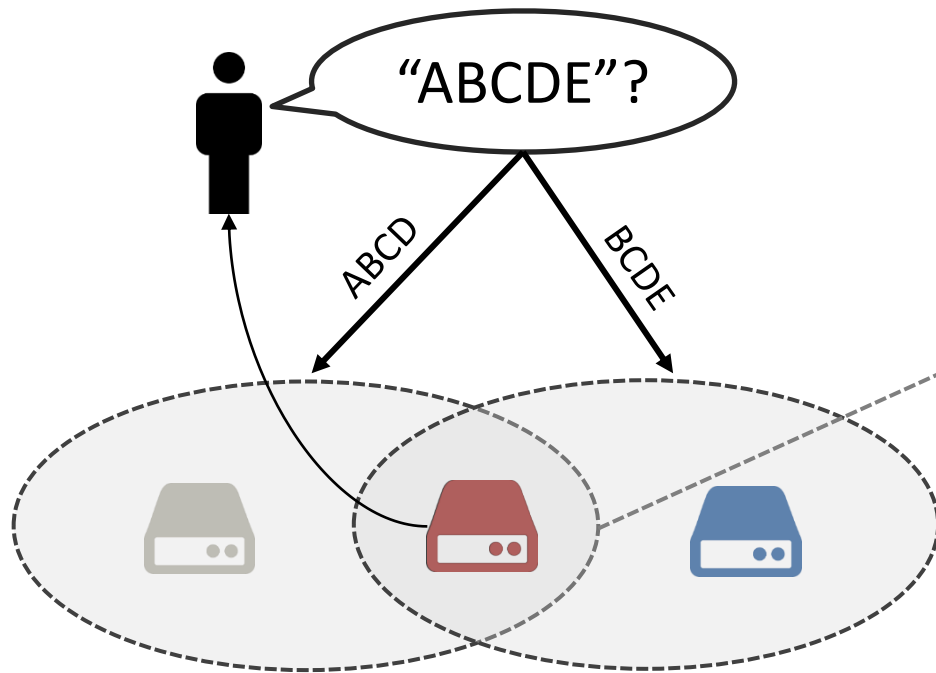HELMHOLTZ CENTER FOR
INFORMATION SECURITY

**YARA rule**

"a.exe" and "exit"
or
"\\x10tes[0-9]" and pe.dll

⬇

**YarIx Search Terms**

or
├─ and
│   ├─ and
│   │   ├─ "a.ex"
│   │   └─ ".exe"
│   └─ "exit"
└─ and
    ├─ "\x10tes"
    └─ pe.dll

✅ Plain Strings (“`a.exe`”, “`exit`”, “`pe.dll`”, …)

✅ Hex Strings (“`{ DE AD BE EF }`”, “`{ CA FE FE BA [2-5] BE FF FF FF }`”)

✅ Regular Expressions (“`calc[0-9a-z]+\.exe`”)

✅ `2 of` (“ABCD”, “ABCE”, “ABCF”, “BCDE”)

✅ Condition Logic

“ABCD” `and` “ABCE” ➜ “ABCD” ∩ “ABCE”
“ABCD” `or` “ABCE” ➜ “ABCD” ∪ “ABCE”

- **Search terms using unsupported features are (e.g., the `not` keyword) over-approximated**

- Offset-free



- Delta Encoding for File IDs

  {1000000, 4, 1, 5, 2, 4, 4 }

- Variable-length 7-bit encoding for File IDs

- Grouping

# Summary

Plain Strings (**"a.exe"**, **"exit"**, **"pe.dll"**, …) ✅
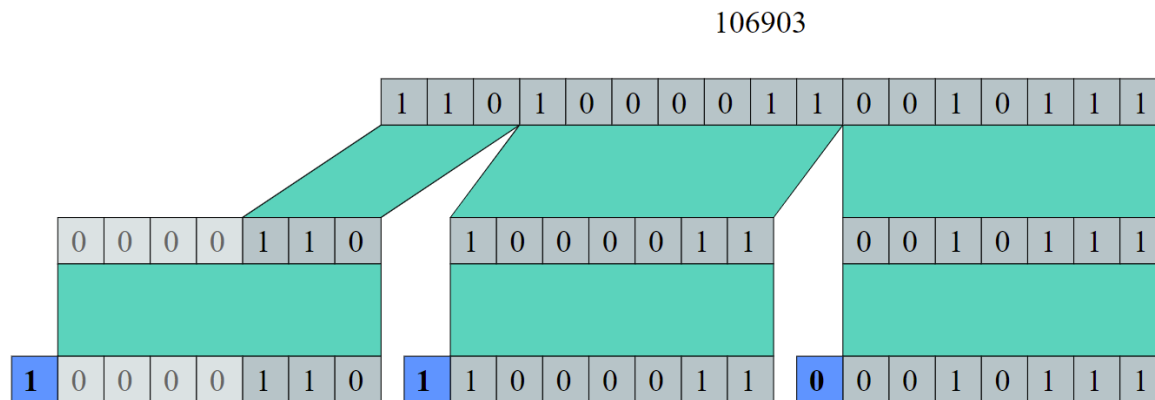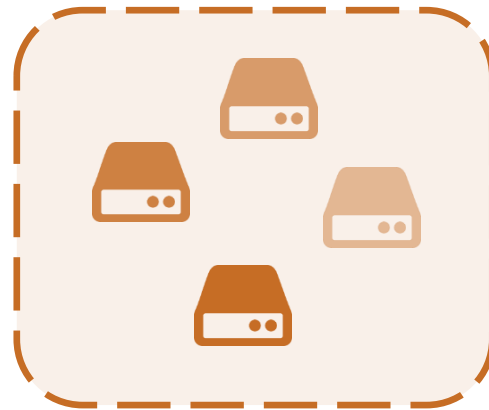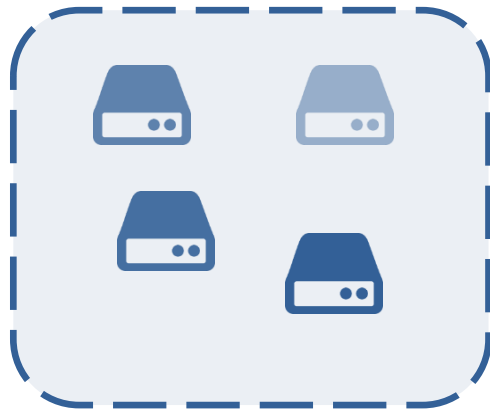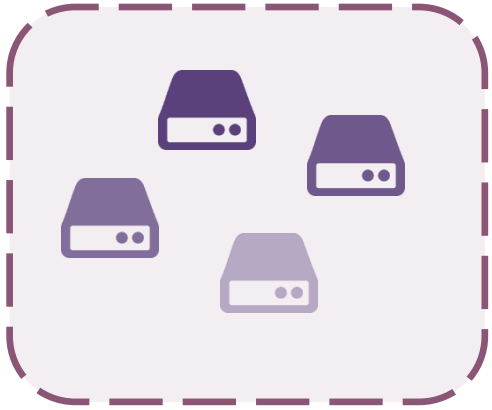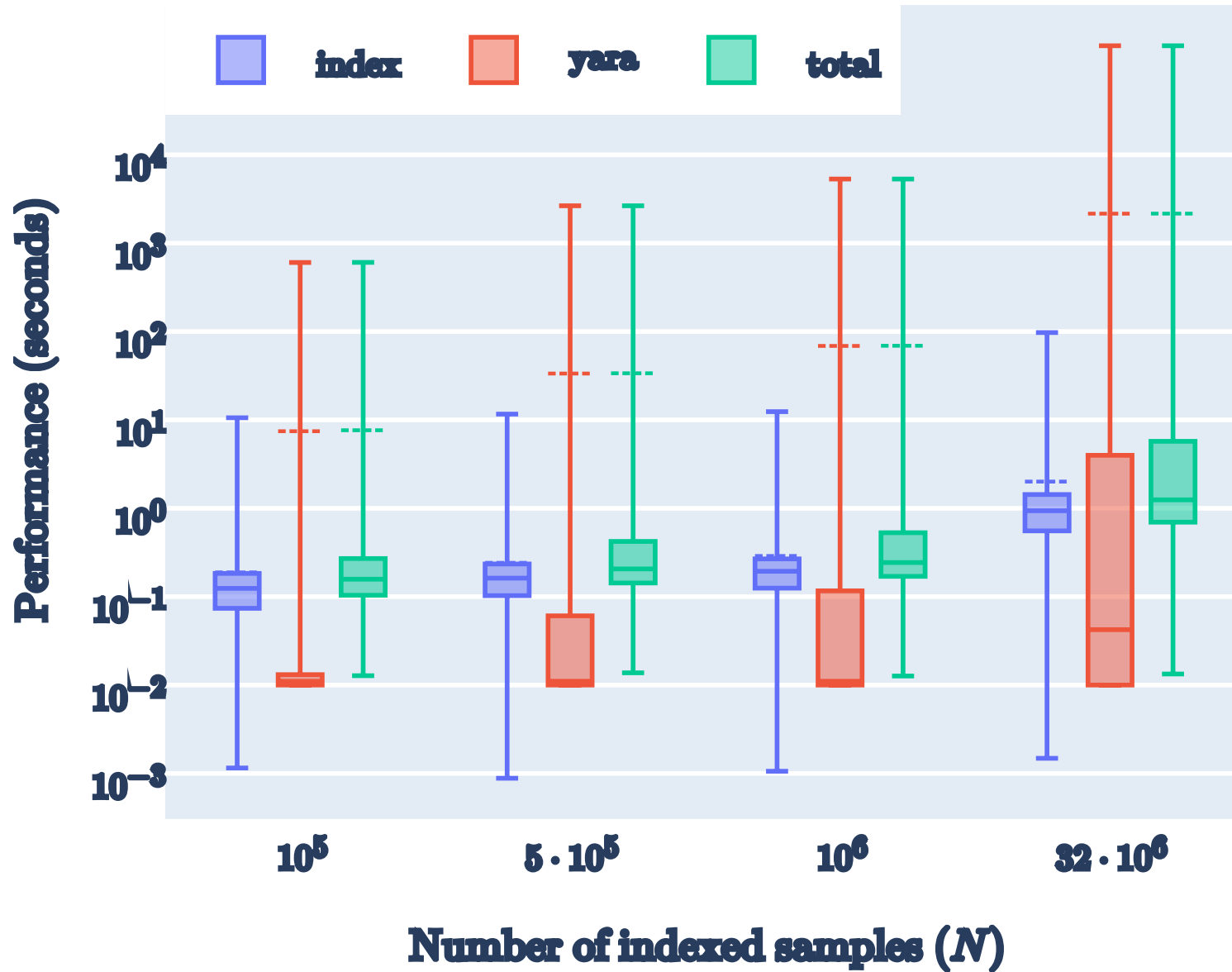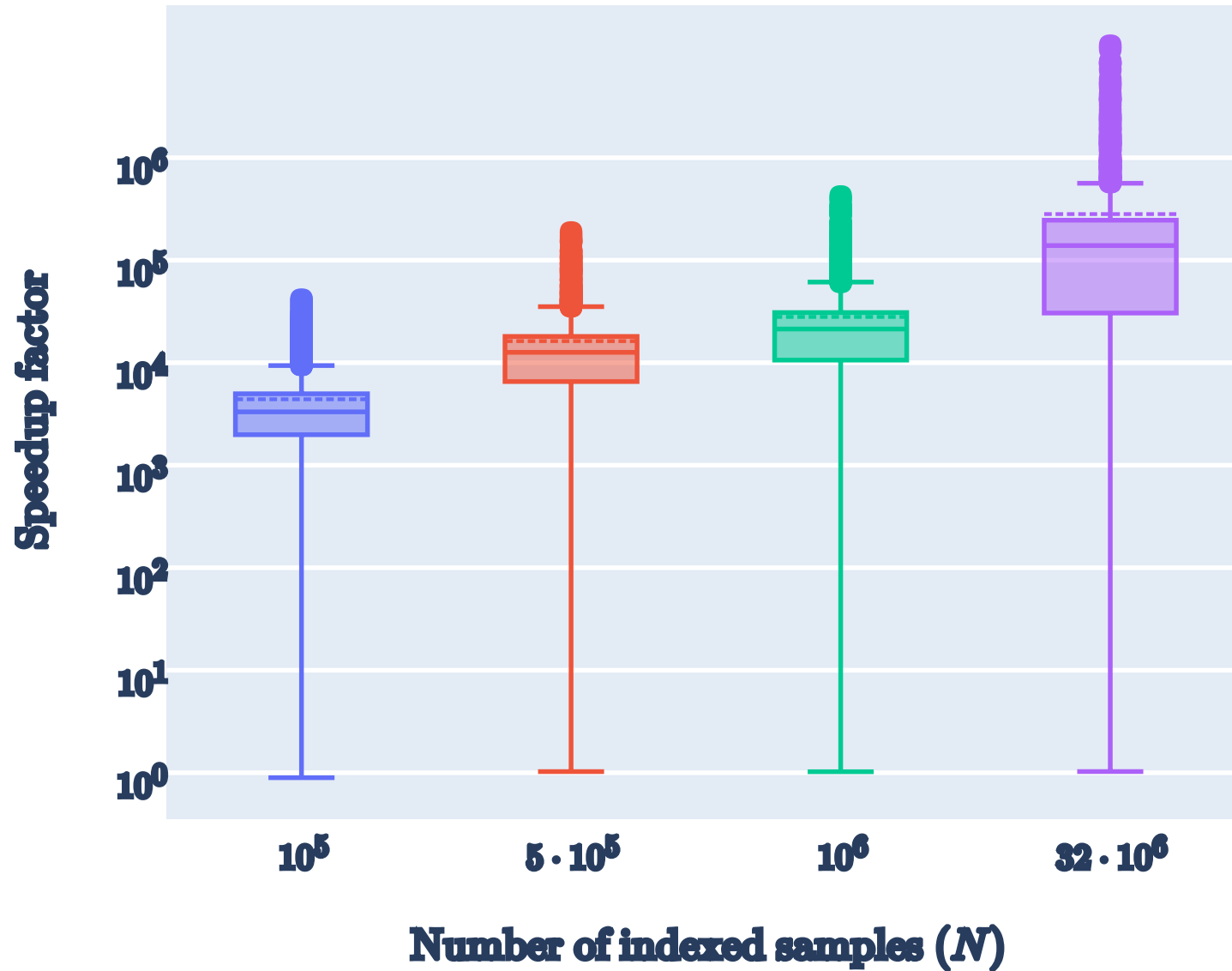
Hex Strings (**"{ DE AD BE EF }"**, **"{ CA FE FE BA [2-5] BE FF FF FF }"**) ✅

Regular Expressions (**"calc[0-9a-z]+\.exe"**) ✅

**2 of ("ABCD", "ABCE", "ABCF", "BCDE")** ✅

Condition Logic ✅



```
            or
          /    \
       and      and
      /   \    /  |  \
   and   "exit" "\x10tes" pe.dll
   /  \
"a.ex" ".exe"
```

| 4-gram | Posting Lists |
|--------|---------------|
| ABCB   |               |
| ABCD   |               |
| ABCE   |               |
| ABCF   |               |
| BCDE   |               |
| CDEA   |               |
| CEGI   |               |
| DEAB   |               |

**https://github.com/mbrengel/yarix**