# THE SCARY SLIDE

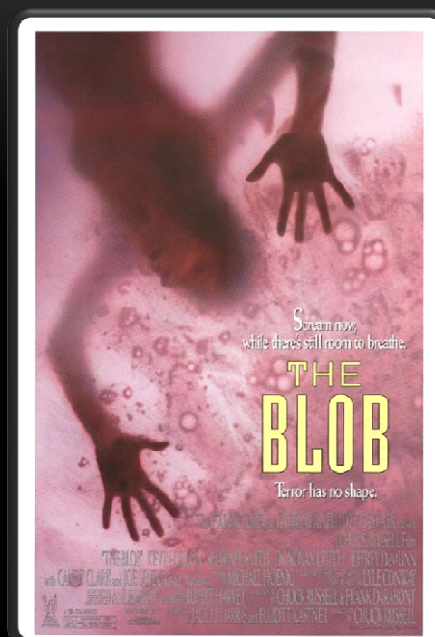- Malware collection size is growing *fast*

    - Analyst teams get  larger but this hinders effective communication: can't double-check with everyone in the team about every new sample (anymore)

    - As a result:

        - Sample assignments are random

        - Signatures become more redundant

        - Malware naming goes downhill, many generic families

# SOLUTION: FAST BINARY BLOB SEARCH

- Think "Bing" for binary (malware) content!

- **FAST** – results in seconds, or less

- **Instant feedback for analysts**

- **MASSIVE** – terabytes of data

- Content agnostic

  - Works on dumps, scripts, JPEGs, anything

  - (we index dumps, for now)

- Design parameter: need a minimum of four contiguous bytes to query

*Microsoft*

DEMO

Microsoft

# NAMING

- Don't know the name? Just grab interesting bytes from the dump and query

# SAMPLE COLLECTION

- Got some kind of marker, not good enough for a detection

  - "{adif}" == 0x7b 0x61 0x64 0x69 0x66 0x7d

- Get samples, write better signature

- Easy to check future samples

*Microsoft*

# INSTA-SIGNATURE TESTING

- Allows instant feedback on refinement of existing signatures
  - Prevent false positives: by instantly checking if patterns match common clean files
    - Works on dumps; **helps prevent in-mem FPs**
  - Limit splash damage: patterns match samples in other families
  - Verify tens of thousands of expected matches in matters of seconds:
    - Much faster than authoring signature, compiling, and then waiting for the product to scan at 50 files/second
- Analyst sig refinement algorithm:
  - $x$ = short fragment, $n$ = hit count
  - while( hitcount( $x$ + extra ) < $n$ )
    - refine(extra);

*Microsoft*

# FAQ

- Q: Can I do boolean combinations?
  - **Yes!** (be careful with **NOT**)
- Q: Is it possible to have wildcards in queries?
  - **Sort of!** Need $\geq$ 4 non-wild contiguous bytes either side.
    - 0x1122334455xxxx6677889900 → 0x1122334455 && 0x6677889900
    - (Under the hood → 0x11223344 && 0x22334455 && 0x66778899 && 0x77889900)
- Q: Can I make case insensitive queries?
  - **Yes…** (sort of, bindex only knows about bytes, but can be done in query frontend)
  - Create all the permutations of upper and lower case bytes, submit the queries and aggregate the results
  - Need $2^x$ queries, where $x$ is the length of the string (and $x \geq 4$)

*Microsoft*

8

# IMPLEMENTATION CHALLENGES

- What index data structures are appropriate for binary content

- How to deal with junk (e.g., compressed or packed data)

- How to make indexing/querying **FAST**

- How to make indexing/querying **SCALE**

- How to deal with the **VOLUME** of data

*Microsoft*

# INDEXING 101

- Each document is assigned a unique Document ID

- This identifier, and the document it refers to, is stored in the Document Map

- Each document is tokenised, and each token is then associated with the document ID in a structure that makes it very efficient to look up tokens. This process is called "**inverting the index**"

Document ID: 512

Document Map

| 512 | \\magrathea\dumps\31\be\31bee721b5c6eab885b0191b8b0c219d736dda72 |
| 513 | \\magrathea\dumps\23\57\235793a658472d4946c833f32867b0749109630e |

*...etc...*

*Microsoft*

# BINDEX APPROACH SIMILAR TO INDEXING ASIAN LANGUAGES

- Asian text typically contains no spaces with which to tokenise
- $n$-gram approach is typical
- We use 4-grams

<div align="center">巨大的云</div>

```
e8  5d  27  59  84  76  91  4e
e8  5d  27
    5d  27  59
        27  59  84
```

*Microsoft*

# SEARCH QUIRKS

- For these two "documents"

  0x11 2233445566778899
  0x1122334400 22334455

- The "query" 0x1122334455 is implemented under the hood as

  0x11223344 && 0x22334455

  so will return a match for both "documents"

- In practice, this doesn't seem to matter

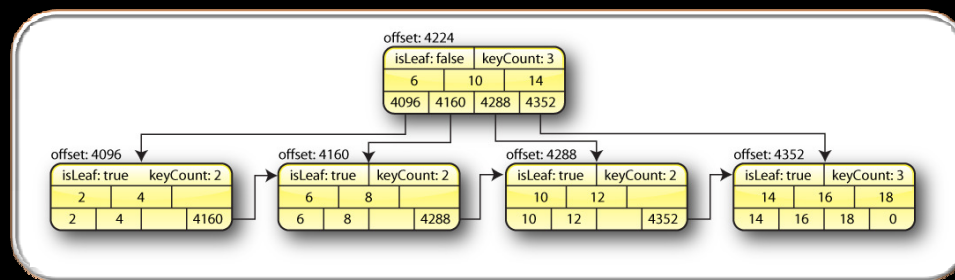  (mostly --- have had FPs with ASCII strings encoded as UTF-16)

  "cloud"
  63 00 6c 00 6f 00 75 00 64 00

  Any doc with UTF-16
  "cl", "lo", "ou" and "ud"
  *Not necessarily adjacent!*

*Microsoft*

# STORING THE N-GRAMS

- B+ trees are ideal
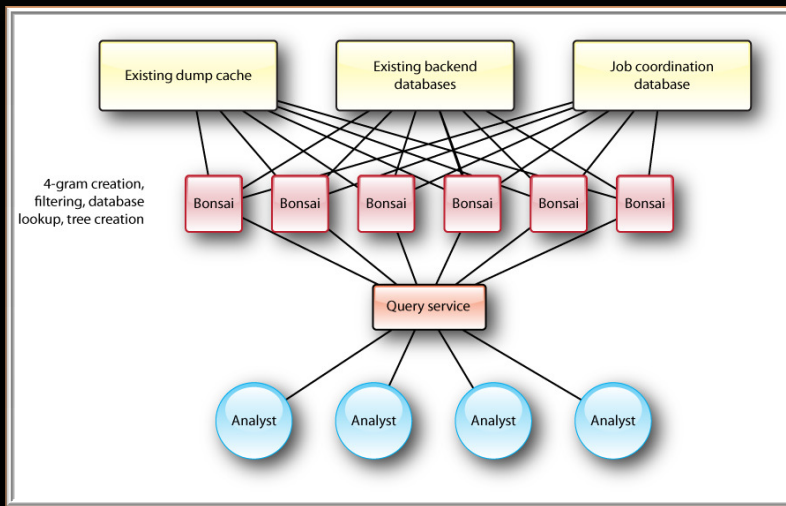
  - Designed for external storage use, very high fan-out minimises disk hits, node size aligned with disk block size

  - Leaf traversal capability helps with duplicate key / value pairs

  - If Document ID is ulong / uint, can jam it directly into the leaves

  - Some implementations keep a parent offset as part of the node (convenient, but bad idea)



*Microsoft*

# SCALABLE ARCHITECTURE

# AVOIDING JUNK

- Don't want to index compressed or encrypted data
    - Useless for search
    - Pollutes and bloats index
- Most malware today is compressed and encrypted, bad to index
- **Solution: we index memory dumps** (existing system already in place generating these)
    - Problem 1: Memory dumps are big (avg 800K, some >20MB)
    - Problem 2: Multiple dumps per file
- Uniquify n-grams per dump (result: BIG reduction in size, ~50% for memory images)
    - Ratio of original n-gram count to unique count good indicator of compression / encryption
    - We tried throwing out n-grams we see more than three times
        - Stops indexing 0x00000000, only makes a small different to count though (1-2%)
    - Multiple dumps can simply be concatenated, but this leads to "linear combination" potential FPs. Not sure if this is a problem in practice.

*Microsoft*

# MANAGING DISK SPACE

- We process new dumps incrementally, and build fixed size trees (bonsai)

    - Easy for distributed construction and search

    - Permits recycling strategy: can simply delete old trees (this is extra good --- delete in B+ trees is so tricky that most books leave it as an "exercise to the reader")

    - Easy to manage disk space, but can only query against the last few days

- We roll over trees once they hit a certain size (say 8G)

    - Enables trees to be constructed entirely in memory, then serialised, for **massive** speedup.

*Microsoft*

# BUILDING THE INDEX

- Tried many, many approaches to building the B+ trees fast. Have to keep up with dump infrastructure

  - First tried SQL Full Text Engine [findex] - fast indexing initially, good query speed but excessive space usage (x15 or more). Seems to break down when things get really large.

  - Using existing databases SQLite, SQL Server vs. rolling our own (roll your own, of course, code reuse is a sin ☺)

  - Native code vs. managed code (complex issue, we tried both, changed mind many times, will let you know my opinion on the day)

  - 64-bit vs. 32-bit trees. This decision affects utterly everything, and we kept picking the wrong one! Solution is a bit subtle…

  - Caching strategies

*Microsoft*

# PERFORMANCE RESULTS

- Around 270K inserts / sec / core on older hardware (>1M on my spiffy i7), but can build many trees at once

  - Approximately one average 700K dump / sec / core

  - (on average, one 700K dump yielded 260K unique 4-grams)

- On average, index is 4x size of data, when indexing dumps using 4-grams

- Scales linearly with cores, provided enough memory

  - Disk not heavily loaded during tree construction

*Microsoft*

# SUMMARY

- BINDEX is
  - practical
  - scalable
  - performant
- BINDEX has a valuable role to play in modern analyst workflow
  - naming
  - sample gathering
  - signature refinement
  - FP testing (including in-memory FPs)
- Maybe even automation…

*Microsoft*

THE END

Microsoft