# DFRWS
## DIGITAL FORENSIC RESEARCH CONFERENCE

Privacy-Preserving Email Forensics

*By*

**Frederik Armknecht, Andreas Dewald and Michael Gruhn**

*Presented At*

The Digital Forensic Research Conference

**DFRWS 2015 USA** Philadelphia, PA (Aug 9th - 13th)

# Privacy-Preserving Email Forensics (PPEF)

Frederik Armknecht (University of Mannheim), Andreas Dewald (Friedrich-Alexander University)

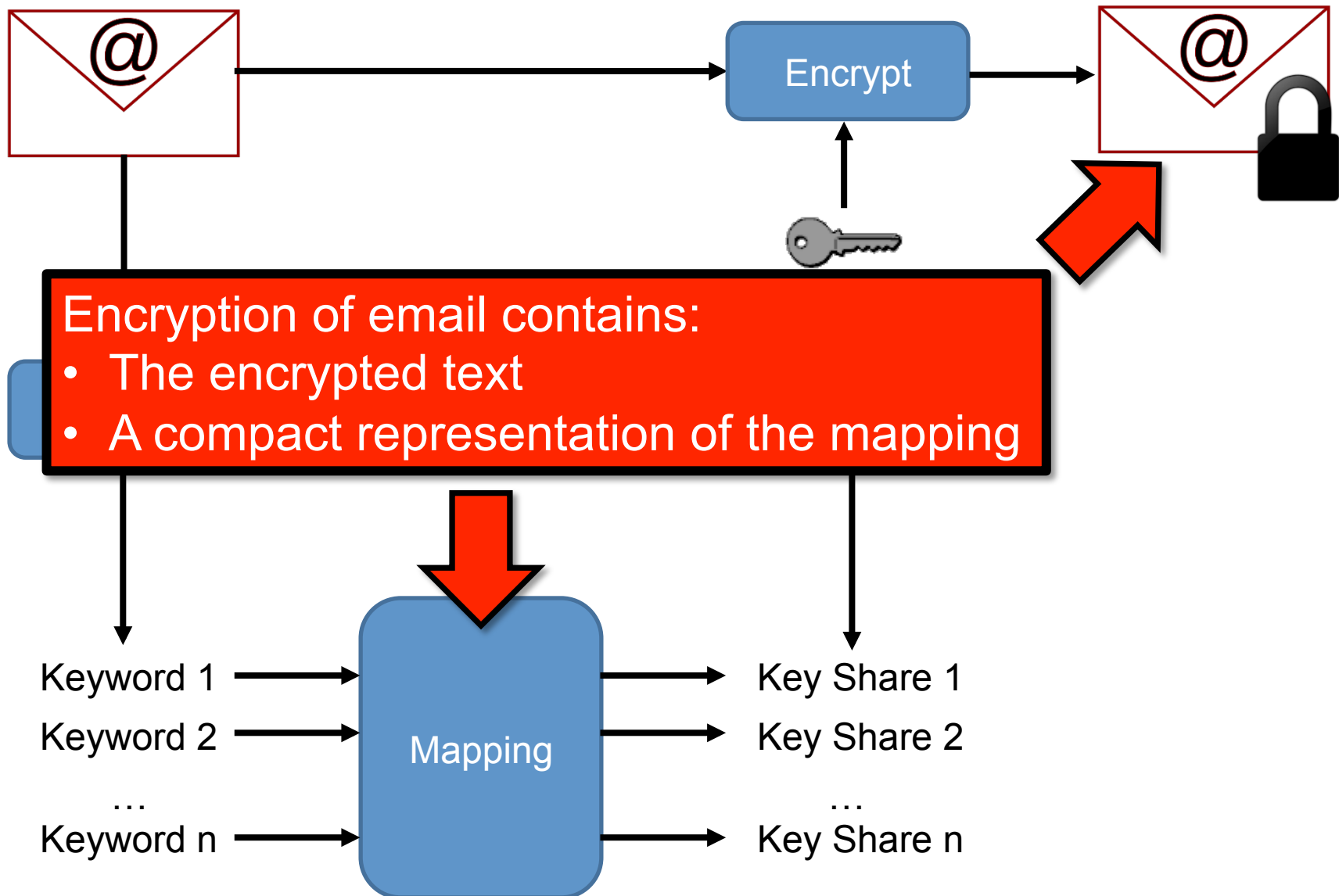**Sven Kälber (Speaker, Friedrich-Alexander University)**

# Agenda

- Idea, motivation & contributions

- The big picture / overall scheme design

- Own implementation details

  - Protection mechanism

  - Extraction mechanism

- Cryptographic building blocks

- Practical implementation and evaluation

- Summary & conclusion
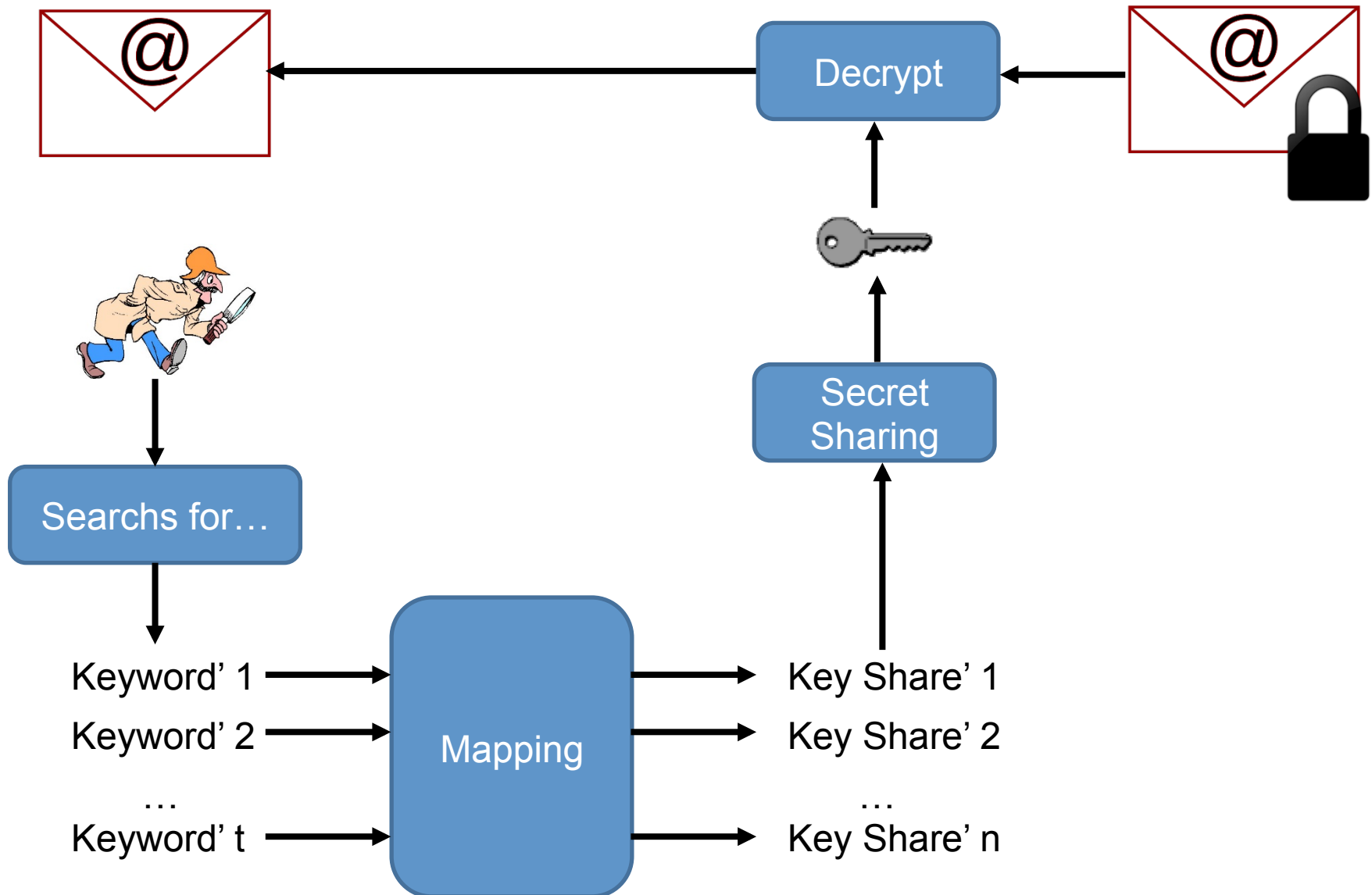
- Limitations & future work

# Idea of PPEF

- Privacy protection of employees in (large-scale) digital forensic investigations

- Revealing of only case relevant information
  - Achieved through strong cryptographic standards

- Operation principle:
  1. Extraction of mailboxes
  2. Encryption of all emails by applying our introduced scheme
  3. Hand over of only encrypted mailboxes to third-party investigators
  4. Decryption of individual emails only possible on $t$ matching keywords

# Motivation

- Private use of corporate e-mail accounts
  - Private e-mails typically contain private and very sensitive data
  - This information is often highly protected by local data protection laws
  - Typically case irrelevant information in private e-mails

- Todays approaches and tools are often limited to filtering, which is not enforced
  - Investigators might read private e-mails by accident or on purpose

- Case „United States v. Carey" (1999)

- Problems of leaving the e-mails at the company's IT
  - Leaks search queries of investigators
  - Is costly and time consuming because of the high degree of interaction needed
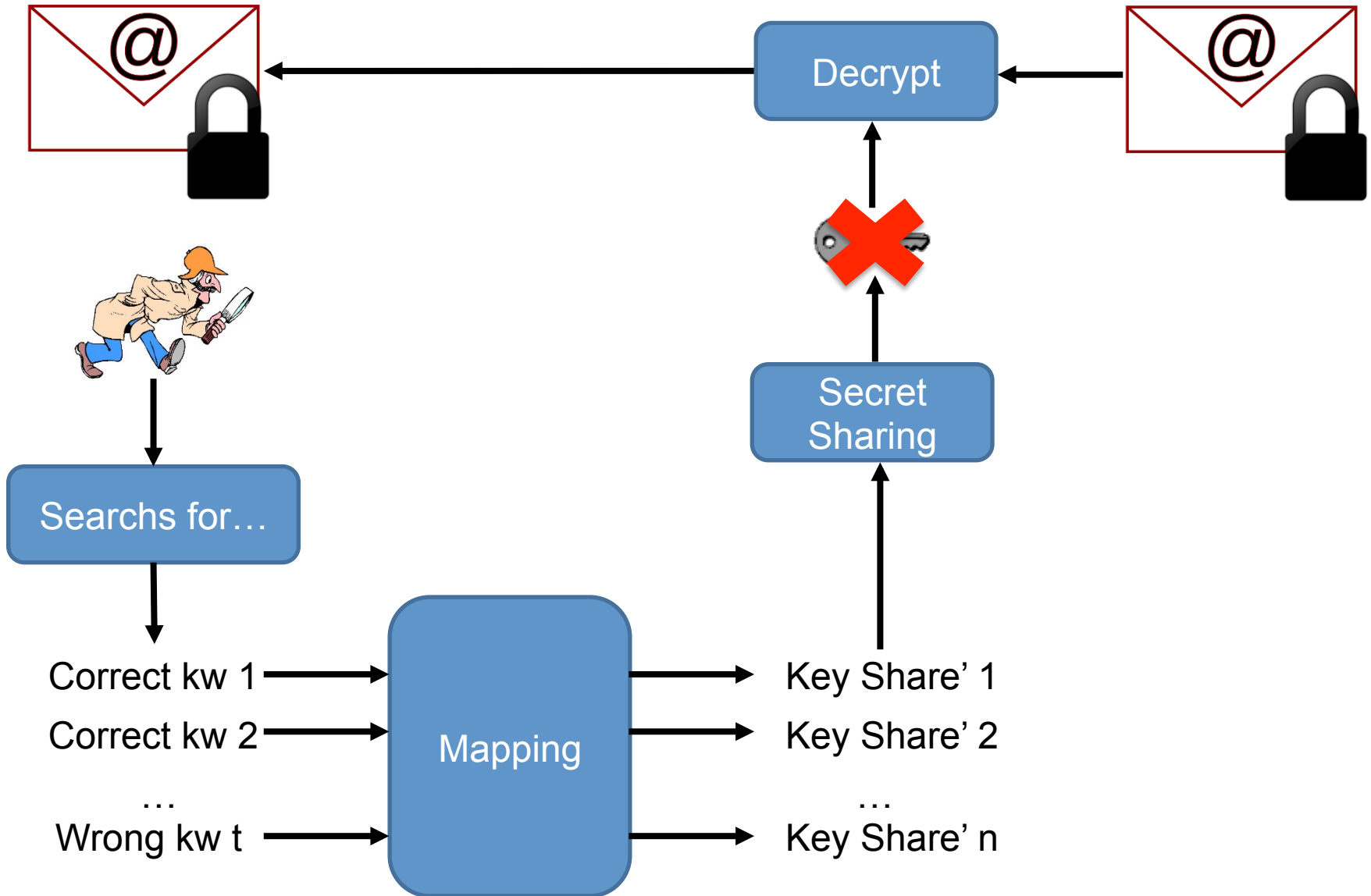  - Trust issues

# Contributions

- Novel approach for privacy-preserving email forensics allowing for **non-interactive** theshold keyword search on **encrypted** emails

- Proof-of-concept implementation in Python and as a Autopsy v3 plug-in

- An evaluation of the practical applicability in terms of:

  - en- / decryption runtime performance
  - introduced storage overhead
  - brute-force / dicitonary attack vulnerability

# Encryption



Encryption of email contains:
- The encrypted text
- A compact representation of the mapping

Keyword 1 → Mapping → Key Share 1
Keyword 2 → Mapping → Key Share 2
… → … 
Keyword n → Mapping → Key Share n

# Decryption success

# Decryption fail

# Details

- Encryption of e-mails (protection mechanism):
  - Each e-mail plaintext $P$ is encrypted to a cyphertext $C$ with an individual secret key $k$.
  - $k$ gets split up in shares and might later be reconstructed
  - Support for blacklisting of commonly used words (e.g. „the")
  - Support for whitelisting of investigation keywords (e.g. „fraud")

- Decryption of e-mails (extraction mechanism):
  - Only possible when the e-mail in question contains at least $t$ keywords.
  - Investigator learns nothing about the secret key of other e-mails upon successfully decrypting one e-mail.
  - Investigator learns nothing about the content of the mail if *t-1 or less keywords* match the content of the e-mail.

# Cryptographic Building Blocks

1. Encryption function:

   - AES-128 in CBC mode used for the encryption of individual e-mails
   - Add characteristic padding *p* as the first block to be decrypted *(e.g. [0,...,0])*

2. Shamir's Secret Sharing

   - Used for splitting the secret key *k* into shares
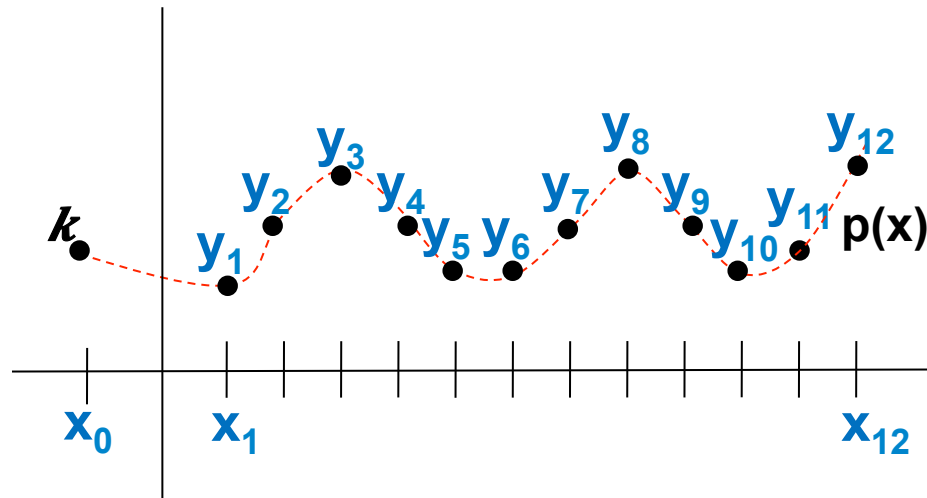   - Details follow on the next slide

3. Mapping

   - Hash function: SHA-256 part of the mapping function
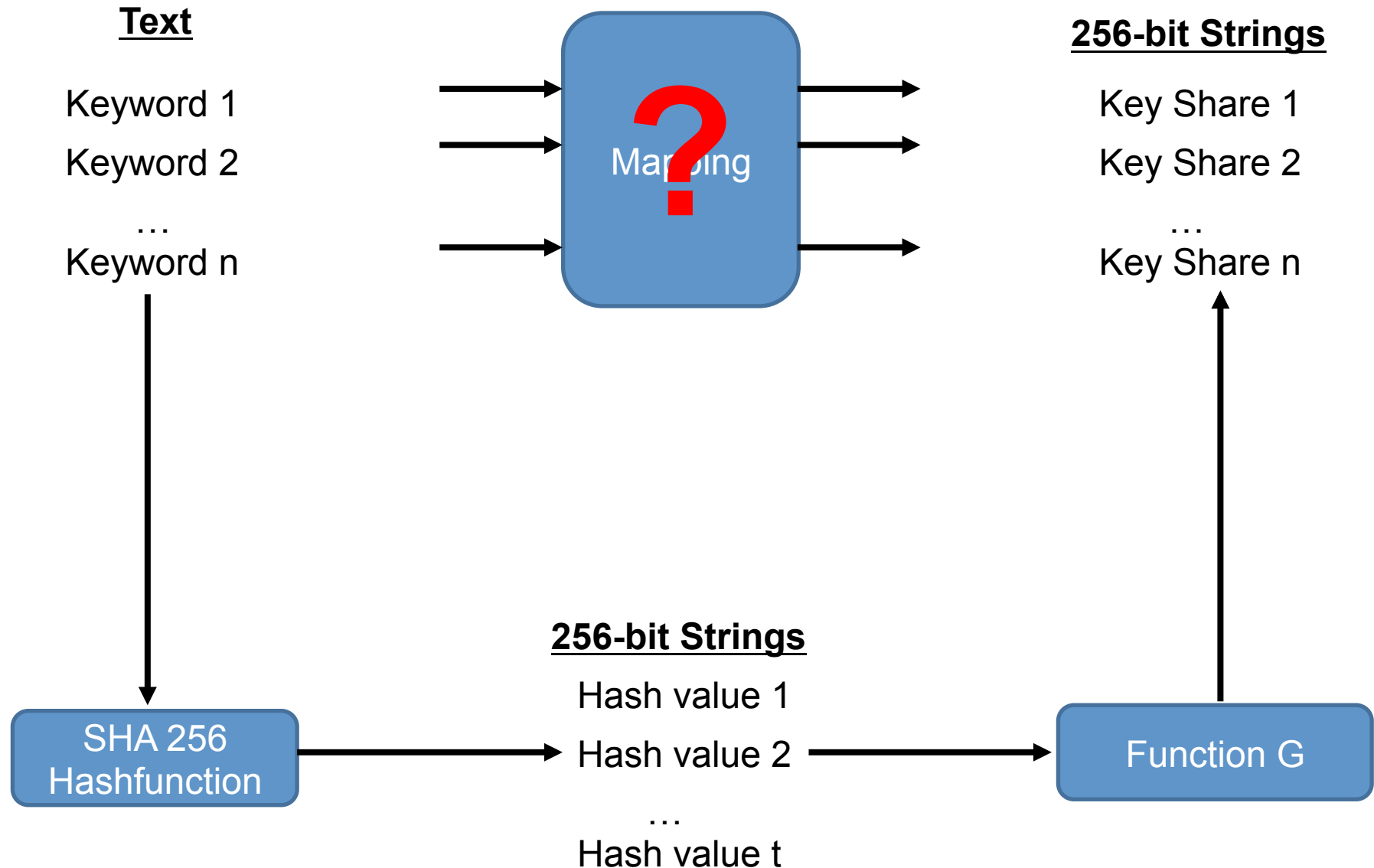   - Further tweaks for efficiency reasons

# Shamir's Secret Sharing

- Functionality:
  - Input: Two integers t ≤ n, a secret k
  - Output: n shares $k_1, ..., k_n$

- Security
  - Given at least t shares, one can reconstruct the secret
  - If less than t shares are known, reconstruction not possible

- Realization
  - Polynomial interpolation of a polynomial of degree t-1

# Working Principle

- Input: Two integers t ≤ n, Secret k

- Choose polynomial $p(x)$ of degree t-1

- Compute shares: $(x_1, y_1),...,(x_n, y_n)$ (here: $n = 12$) with $y_i = p(x_i)$

- Reconstruction from t shares:

  - Interpolate $p(x)$

  - Compute $k = p(x_0)$

# The Mapping Function

**Text**

Keyword 1

Keyword 2

…

Keyword n

Mapping **?**

**256-bit Strings**

Key Share 1

Key Share 2

…

Key Share n

**256-bit Strings**

Hash value 1

Hash value 2

…

Hash value t

SHA 256
Hashfunction

Function G

# Function G

- Task: Map 256-bit hash values to 256-bit shares ($x{\downarrow}i$, $y{\downarrow}i$)

- Approach:
  - Interpret hash values as $(x{\downarrow}i, z{\downarrow}i)$ (128-bit + 128-bit)
  - Use the values $x{\downarrow}i$ to compute shares $(x{\downarrow}i, y{\downarrow}i) = (x{\downarrow}i, p(x{\downarrow}i))$
  - Find mapping $g(x)$ such that $g(x{\downarrow}i) = y{\downarrow}i$ XOR $z{\downarrow}i$
  - Function $G(h{\downarrow}i) = G(x{\downarrow}i, z{\downarrow}i) = (x{\downarrow}i, g(x{\downarrow}i)$ XOR $z{\downarrow}i) \rightarrow (x{\downarrow}i, y{\downarrow}i)$

- Getting the mapping:
  - Core idea: compute polynomial $g(x)$ such that $g(x{\downarrow}i) = y{\downarrow}i$ XOR $z{\downarrow}i$
  - Problem: requires to interpolate polynomial of degree $n$ $\rightarrow$ effort is $O(n{\uparrow}3)$, too slow
  - Idea: Split range of $x$ into $l$ subsets, e.g. determined by the $l$ last bits
  - Interpolate polynomials $g{\downarrow}j(x)$ for each subset
  - Effort: interpolate $l$ polynomials, each of degree $\approx n/l$
  - Overall effort: $l \cdot (n / l){\uparrow}3 = n{\uparrow}3 / l{\uparrow}2$

# Practical implementation

1. Python en- / decrytion of mailboxes

   Supported mailbox formats:

   mbox

   pst

   MH (RFC 822)

   Maildir



2. PPEF plugin for Autopsy v3

# PPEF Autopsy plug-in

# PPEF Autopsy plug-in

# PPEF Autopsy plug-in

# PPEF Autopsy plug-in

# PPEF Autopsy plug-in

# Evaluation

The data set used in our evaluation consists of 5 different mailboxes:

- *Apache* – httpd user mailing list (75724 e-mails)
- *Work* – personal work e-mails (1590 e-mails)
- *A, B, C* – private e-mail accounts (511, 349, 83 e-mails)

## Evaluations:

1. Encryption runtime performance

2. Encryption storage overhead

3. Search / decryption runtime performance

4. Brute-force attack performance

# Encryption Performance

- Time (in seconds) to encrypt the corresponding emails of each account

- Average encryption rate: 13.5 emails/sec
  - Encryption of large mailboxes might take several hours (< 2h for 75724 e-mails) , but only needs to be done once!

|     | Apache [s] | Work | A | B | C |
| --- | --- | --- | --- | --- | --- |
| Min | 0.004 | 0.005 | 0.005 | 0.005 | 0.005 |
| Max | 31.745 | 1.403 | 1.932 | 1.117 | 0.460 |
| Avg | 0.082 | 0.136 | 0.122 | 0.110 | 0.173 |
| Med | 0.072 | 0.115 | 0.101 | 0.074 | 0.150 |
| $\sigma$ | 0.133 | 0.120 | 0.132 | 0.153 | 0.071 |
| $\Sigma$ | 6243.511 | 217.242 | 62.842 | 38.535 | 14.367 |

# Encryption storage overhead

- Encryption with AES does not add much storage overhead (33 – 48 bytes per mail)

- Main storage overhead factor is the mapping function (on average 582.4 bytes per mail)

- Average storage overhead: 5.2 %

|                | Apache   | Work     | A      | B      | C     |
|----------------|----------|----------|--------|--------|-------|
| Size Raw [KB]  | 376, 551 | 418, 680 | 16, 386 | 47, 486 | 6, 676 |
| Size PPEF [KB] | 418, 870 | 420, 418 | 16, 885 | 47, 821 | 6, 806 |
| Overhead       | 11.2 %   | 0.4 %    | 3.0 %  | 0.7 %  | 1.9 % |

# Search / decryption performance

- Time (in seconds) to search each e-mail for 3 keywords and decrypt matching e-mails

- Average search and decryption rate: 98 mails/sec
  - Searches on large mailboxes take time (< 15min) but are still feasible

|  | Apache [s] | Work | A | B | C |
|---|---|---|---|---|---|
| Min | 0.0090 | 0.0096 | 0.0098 | 0.0097 | 0.0098 |
| Max | 0.0598 | 0.1645 | 0.0139 | 0.1508 | 0.0148 |
| Avg | 0.0115 | 0.0137 | 0.0114 | 0.0123 | 0.0117 |
| Med | 0.0115 | 0.0117 | 0.0113 | 0.0113 | 0.0116 |
| $\sigma$ | 0.0007 | 0.0103 | 0.0007 | 0.0086 | 0.0009 |
| $\Sigma$ | 876.8591 | 21.7977 | 5.8650 | 4.2982 | 0.9750 |

# Attack performance evaluation

- Brute-force attacks to decrypt the whole mailbox ($\pi$=0.99) or a random half of the mailbox ($\pi$=0.5)

- Using 4 different vocabularies
  - Oxford English Dictionary (171,476 words)
  - 50 % of the Oxford English Dictionary (85,738 words)
  - Vocabulary in daily speech edu. person (20,000 words)
  - Vocabulary of uneducated person (10,000 words)

| $\pi$ | $N$ | Apache | Work | A | B | C |
|---|---|---|---|---|---|---|
| 0.99 | 171,476 | $1.15 \cdot 10^8$ | $3.26 \cdot 10^5$ | $1.23 \cdot 10^5$ | $1.17 \cdot 10^5$ | 5,373.15 |
| 0.50 | 171,476 | $1.73 \cdot 10^7$ | 49,072.84 | 18,565.34 | 17,638.94 | 808.74 |
| 0.99 | 85,738 | $1.44 \cdot 10^7$ | 40,753.36 | 15,418.00 | 14,648.51 | 671.63 |
| 0.50 | 85,738 | $2.17 \cdot 10^6$ | 6,133.99 | 2,320.64 | 2,204.82 | 101.09 |
| 0.99 | 20,000 | $1.83 \cdot 10^5$ | 517.23 | 195.68 | 185.91 | 8.52 |
| 0.50 | 20,000 | 27,510.37 | 77.85 | 29.45 | 27.98 | 1.28 |
| 0.99 | 10,000 | 22,843.44 | 64.64 | 24.46 | 23.24 | 1.07 |
| 0.50 | 10,000 | 3,438.28 | 9.73 | 3.68 | 3.50 | 0.16 |

# Summary / Conclusion

- We proposed a novel approach for privacy-preserving email forensics allowing for **non-interactive** theshold keyword search on **encrypted** e-mails.

- We developed a **prototype-implementation** in Python and an **Autopsy plug-in** that supports multiple well-known mailbox formats.

- We **evaluated the practical applicability** in terms of en- / decryption performance, storage overhead and brute-force vulnerability.

  - Sufficiently large mailboxes are well protected against dictionary (brute-force) attacks

# Limitations / Future work

Limitations:

- Scheme based on keyword searches, therefore prone to spelling errors

- No wildcard operator or regular expression possible that allows for more advanced search queries

- Brute-force / dictionary attacks possible


Future work:

- Support for wildcard usage within the search keywords

# Thank you for your attention!

Speaker:

Sven Kälber - [sven.kaelber@cs.fau.de](mailto:sven.kaelber@cs.fau.de)

Authors:

Frederik Armknecht – [armknecht@uni-mannheim.de](mailto:armknecht@uni-mannheim.de)

Andreas Dewald – [andreas.dewald@fau.de](mailto:andreas.dewald@fau.de)

# Function G

- Task: Map 256-bit hash values to 256-bit key shares

- Approach
  - Interpret hash values as (128 bit + 128 bit)
  - Use the values to compute key shares
  - Find a mapping such that
  - Function

- Getting the Mapping
  - Core idea: compute polynomial such that
  - Problem: requires to interpolate polynomial of degree n  effort is in , too
  - Idea: Split range of into subsets, e.g., determined by the first bits
  - Interpolate polynomials for each subset
  - Effort: interpolate polynomials, each of degree .
  - Overall effort: