



SSDEEPER: EVALUATING AND IMPROVING SSDEEP

By

Carlo Jakobs, Martin Lambertz and Jan-Niclas Hilgert

carlo.jakobs@fkie.fraunhofer.de, martin.lambertz@fkie.fraunhofer.de, jan-niclas.hilgert@fkie.fraunhofer.de

From the proceedings of

The Digital Forensic Research Workshop

DFRWS 2022 VIRTUAL USA

(July 11th - 14th)

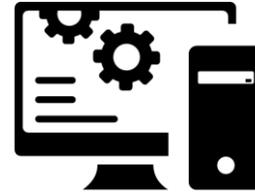
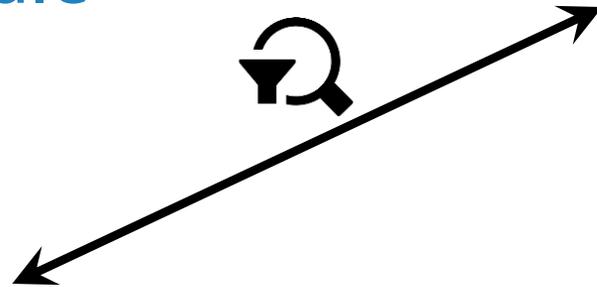
Introduction

The problem of scale



Introduction

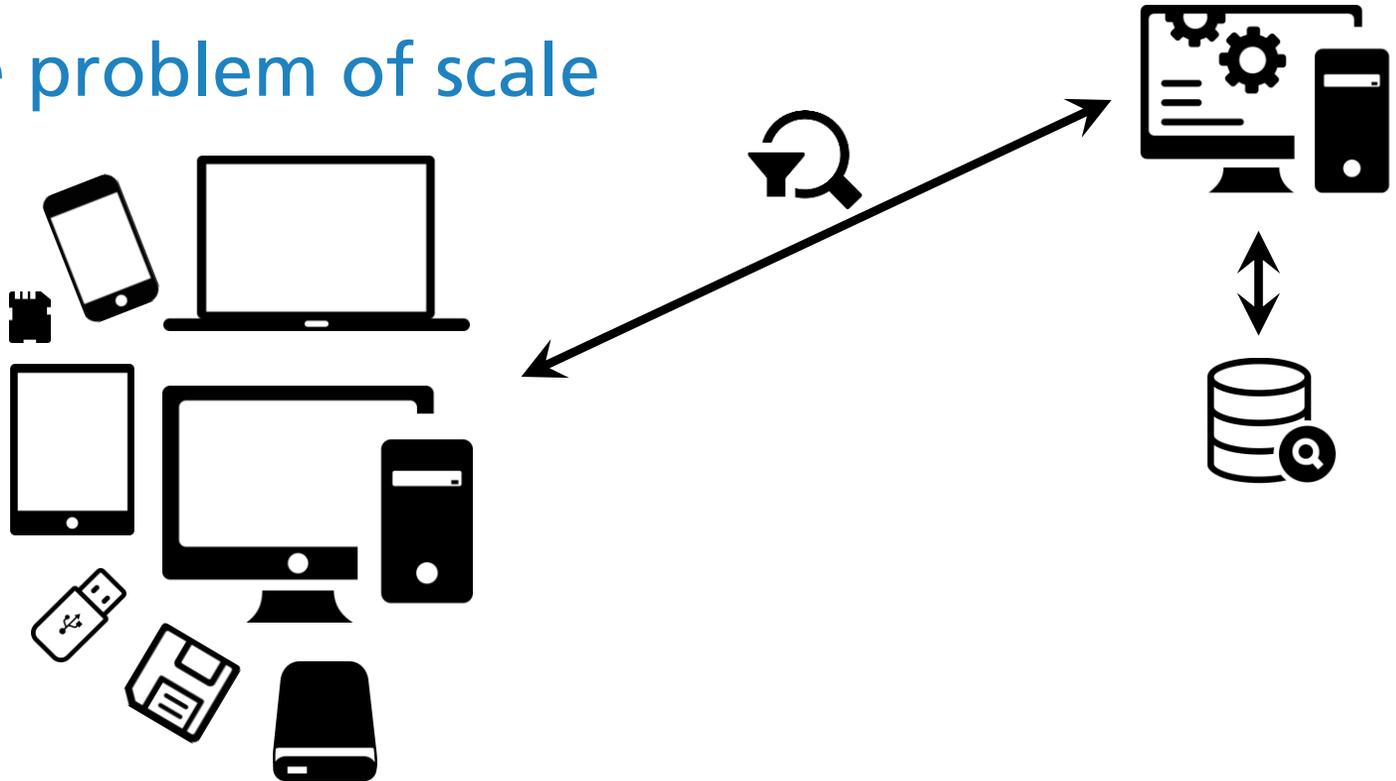
The problem of scale



Automated systems

Introduction

The problem of scale

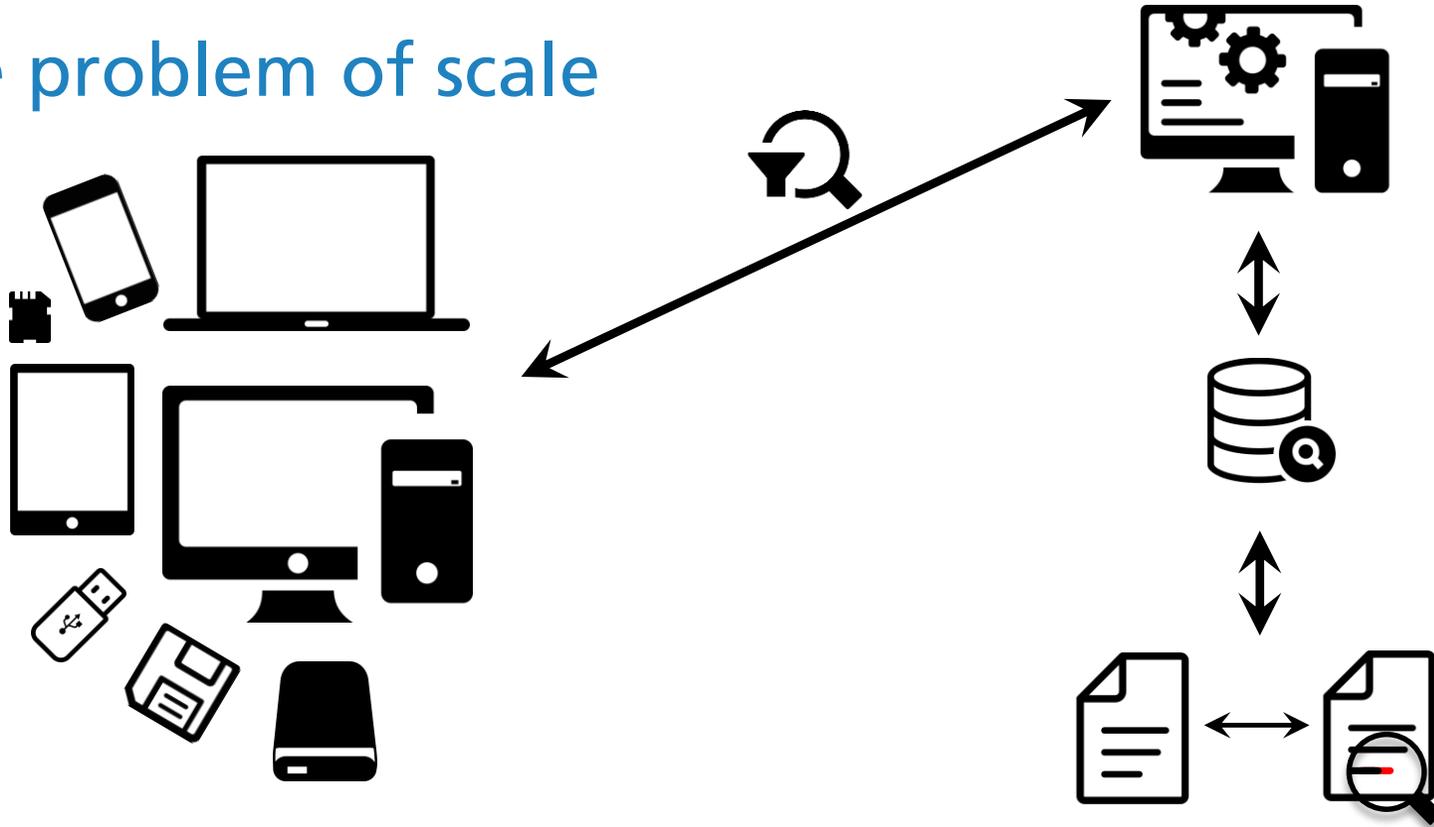


Automated systems

Traditional technique

Introduction

The problem of scale



Automated systems

Traditional technique

Limitation

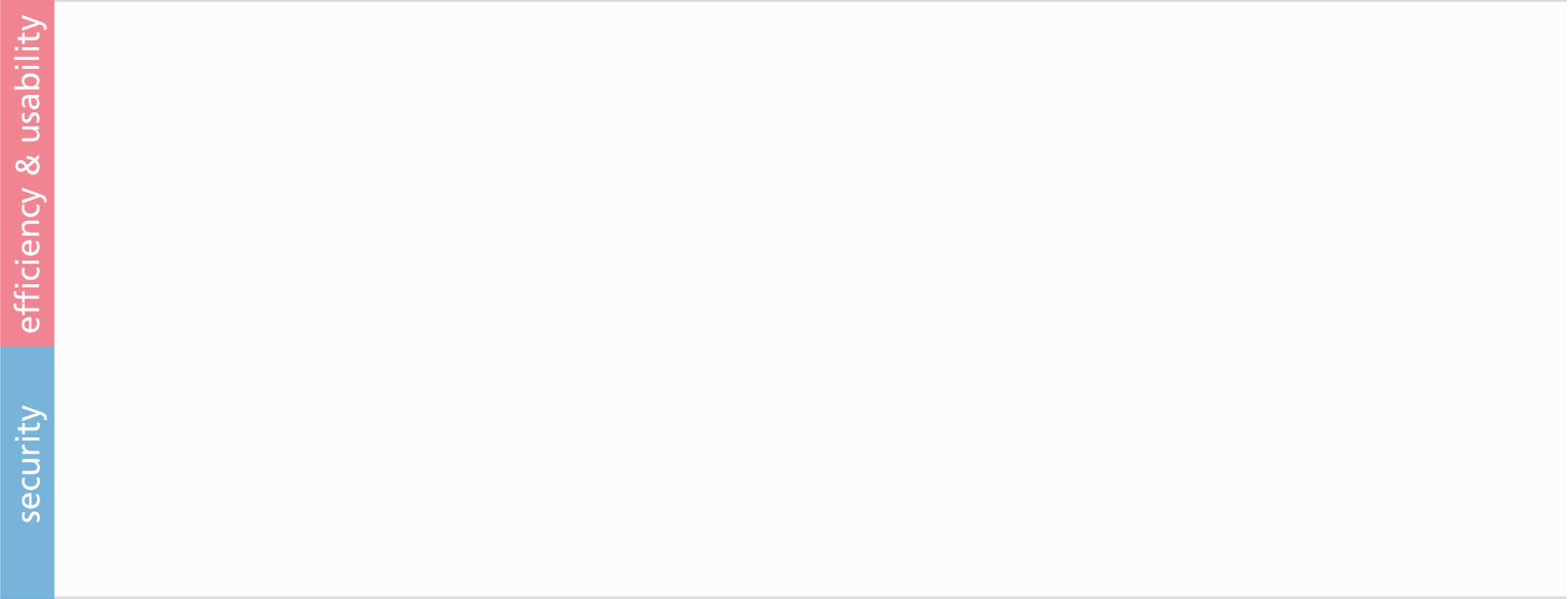
Introduction – Similarity

- What is **similarity**?
 - **Lower level** approximates the similarity between two **byte sequences**
 - **File type agnostic** and **fast**
 - **Higher level** attempts to interpret the **actual content**
 - **File type specific**, requires **parsing** and **interpretation** of bytes

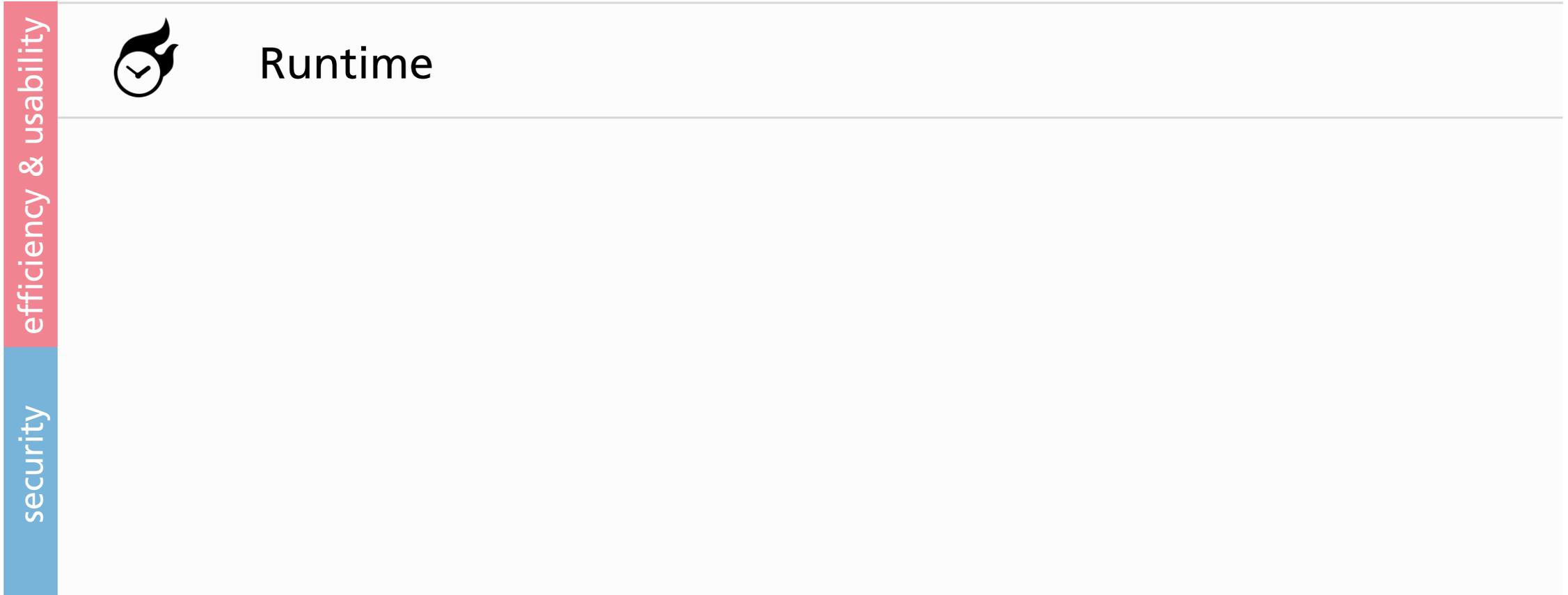
Introduction – Similarity Hashing

- What is **similarity**?
 - **Lower level** approximates the similarity between two **byte sequences**
 - **File type agnostic** and **fast**
 - **Higher level** attempts to interpret the **actual content**
 - **File type specific**, requires **parsing** and **interpretation** of bytes
- **Similarity Hashing**
 - On the **lower level** of abstraction
 - Similarity Preserving Hash Function (**SPHF**)
 - Comparison Function (**CF**)

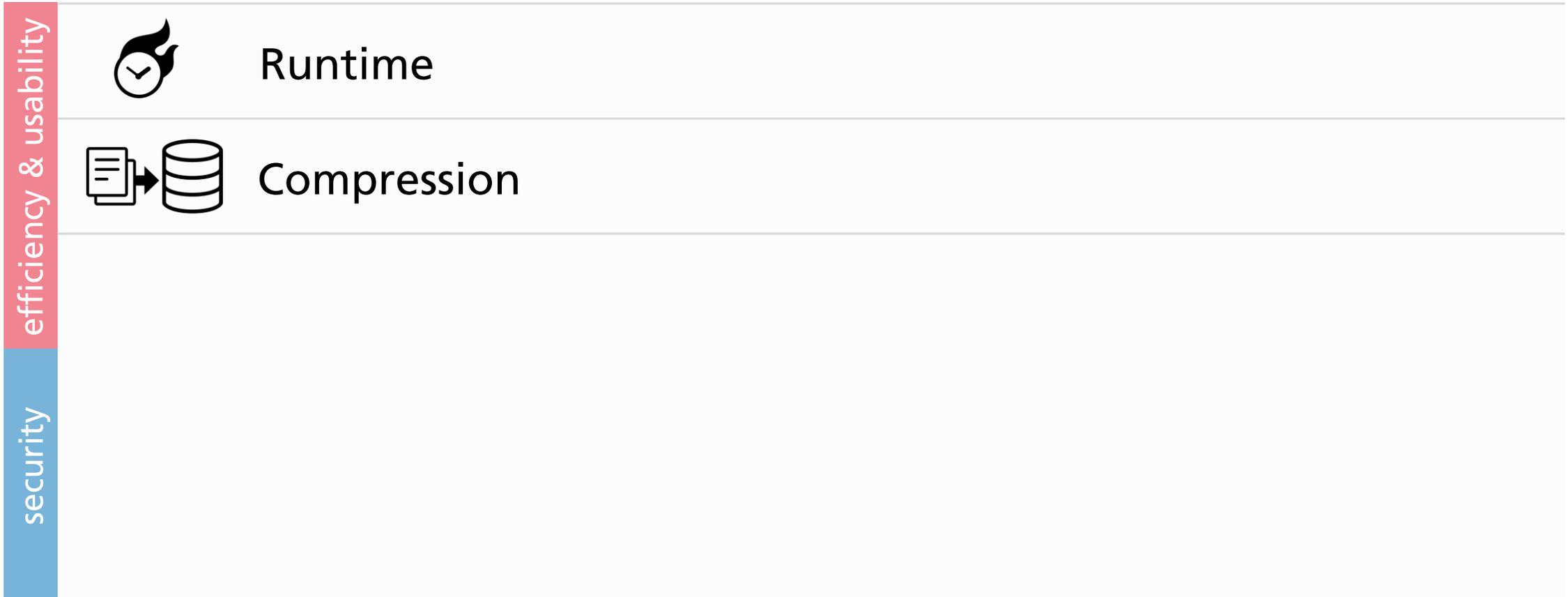
Similarity Hashing – Properties



Similarity Hashing – Properties



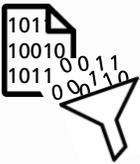
Similarity Hashing – Properties



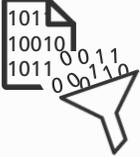
Similarity Hashing – Properties

efficiency & usability		Runtime
		Compression
		Accuracy
security		

Similarity Hashing – Properties

efficiency & usability		Runtime
		Compression
		Accuracy
security		Coverage

Similarity Hashing – Properties

efficiency & usability		Runtime
		Compression
		Accuracy
security		Coverage
		Obfuscation resistance

Similarity Hashing – *ssdeep* – SPHF

window

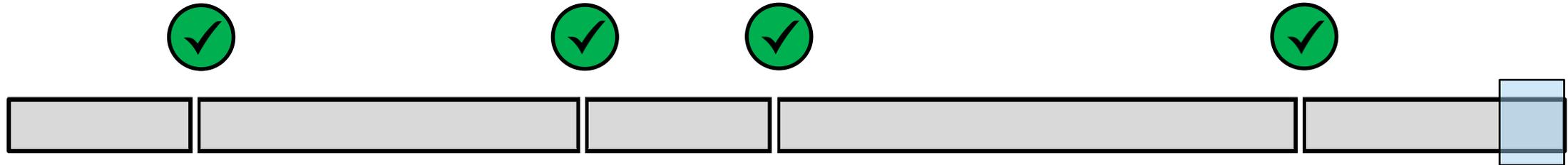


rolling hash function
Adler32 checksum

$$h \equiv -1 \pmod{b}$$

traditional hash function
FNV-1 (6 bits)

Similarity Hashing – *ssdeep* – SPHF



rolling hash function
Adler32 checksum

$$h \equiv -1 \pmod{b}$$

traditional hash function
FNV-1 (6 bits)

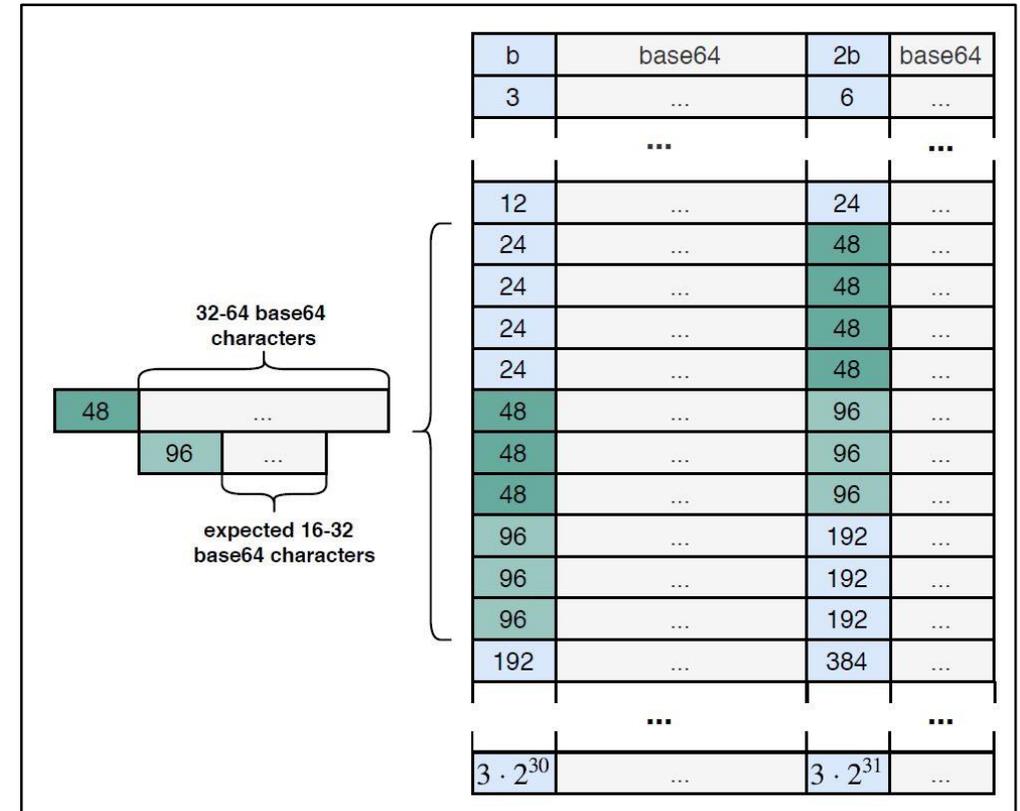
<i>b</i>	t	Q	y	N	F
----------	---	---	---	---	---

- Recalculation if $len < 32$
 - Use $b = b/2$
 - `blocksize:hash:hash,filename`
- Variety of improvements

Similarity Hashing – ssdeep

■ CF

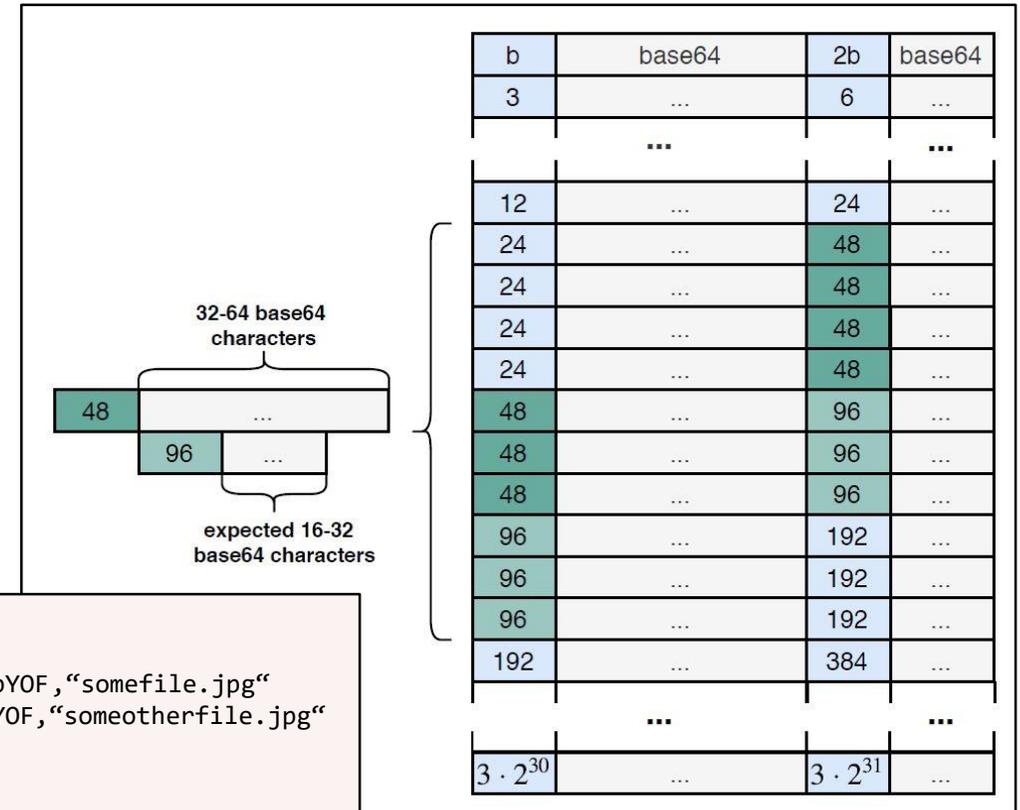
- Matching block sizes
- Levensthein edit distance – string comparison



Similarity Hashing – ssdeep

■ CF

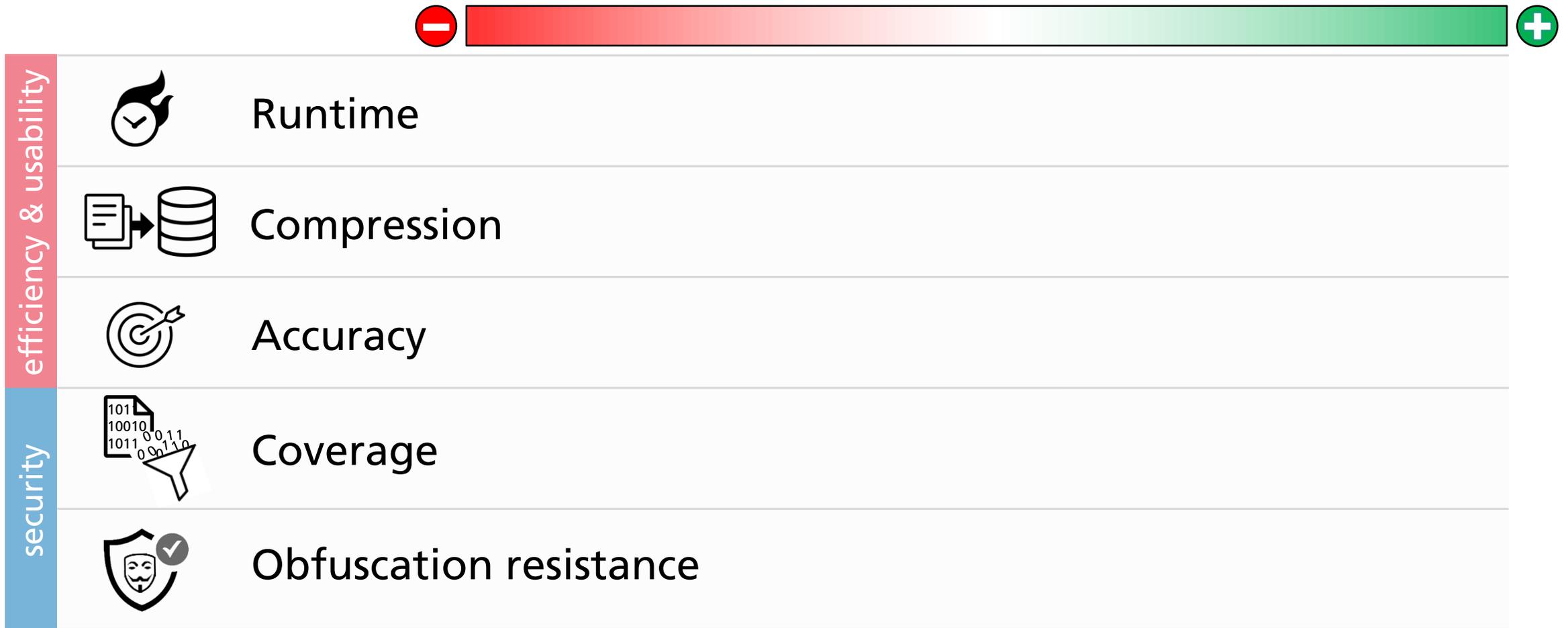
- Matching block sizes
- Levensthein edit distance – string comparison



```
$ ssdeep somefile.jpg someotherfile.jpg
ssdeep,1.1-blocksize:hash:hash,filename
192:tQyNFgeWKVwVyF18TYjEiHRElsOMuKyN46LouR1dPbCNXf4pMt91j3IOYEVF:c1xA1+00Ay1LVXPbAXf7t9ZbYOF,"somefile.jpg"
192:tQyNFgeWKVuyF18TYjEiHRElsOMuKyN46LouR1dPbCNXf4pMt91j3IOYEVF:c1xA1+00Ay1LVXPbAXf7t9ZbYOF,"someotherfile.jpg"

$ ssdeep somefile.jpg someotherfile.jpg > known_hashes.txt
$ ssdeep -m known_hashes.txt someotherfile.jpg
someotherfile.jpg matches somefile.jpg (99)
someotherfile.jpg matches someotherfile.jpg (100)
```

Similarity Hashing – ssdeep – Properties

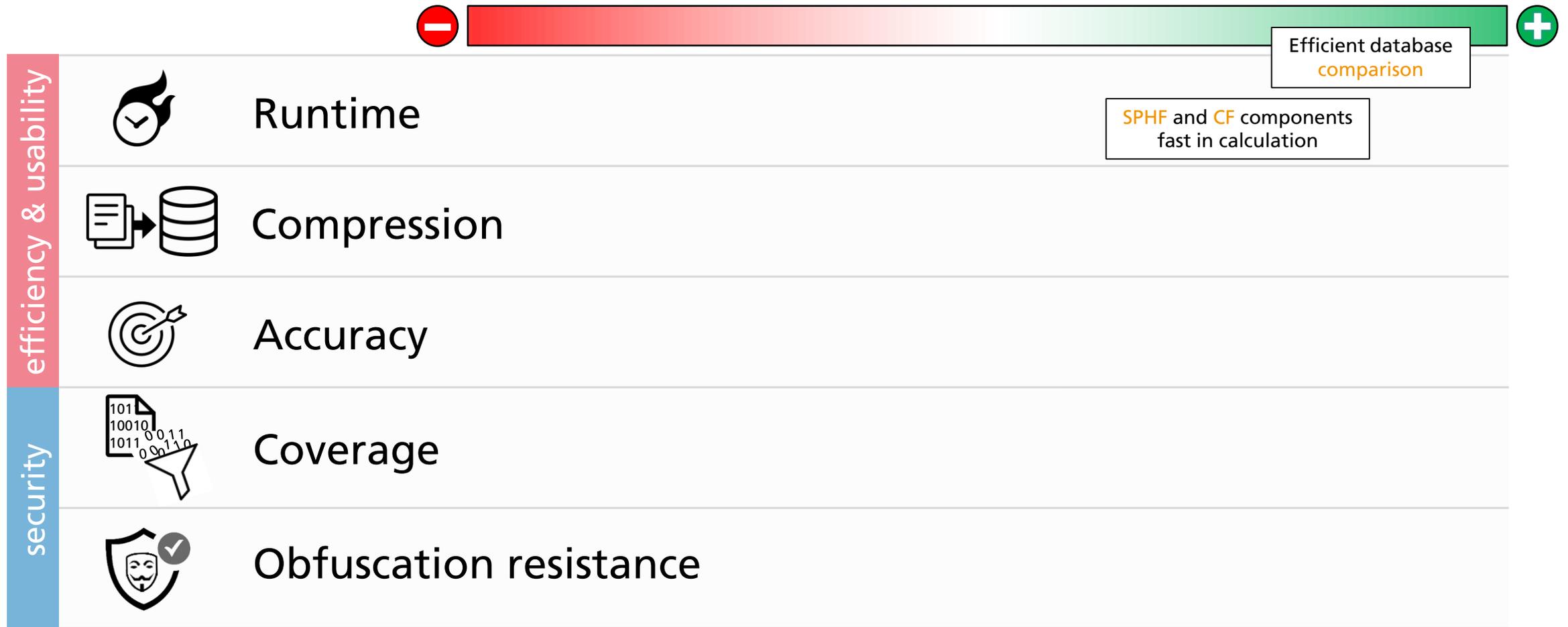


Similarity Hashing – ssdeep – Properties

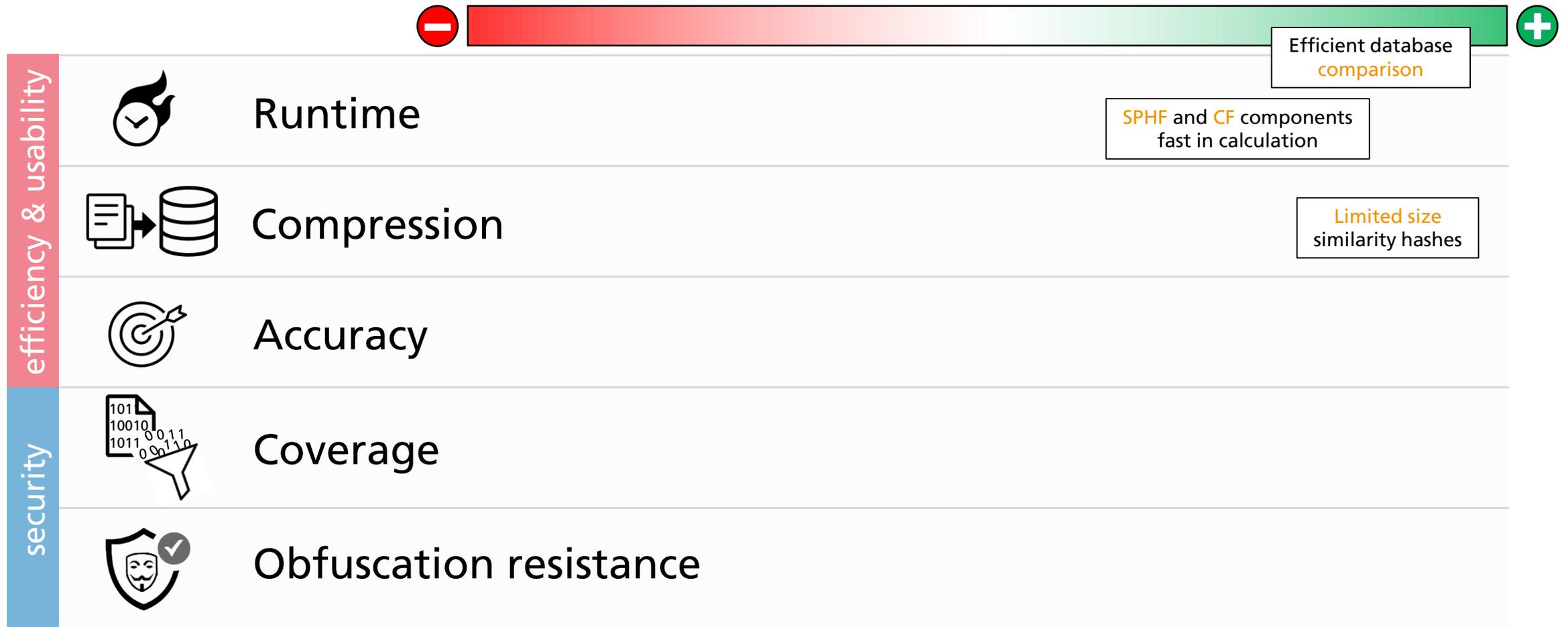


efficiency & usability		Runtime	<div data-bbox="1786 472 2175 562">SPHF and CF components fast in calculation</div>
		Compression	
		Accuracy	
security		Coverage	
		Obfuscation resistance	

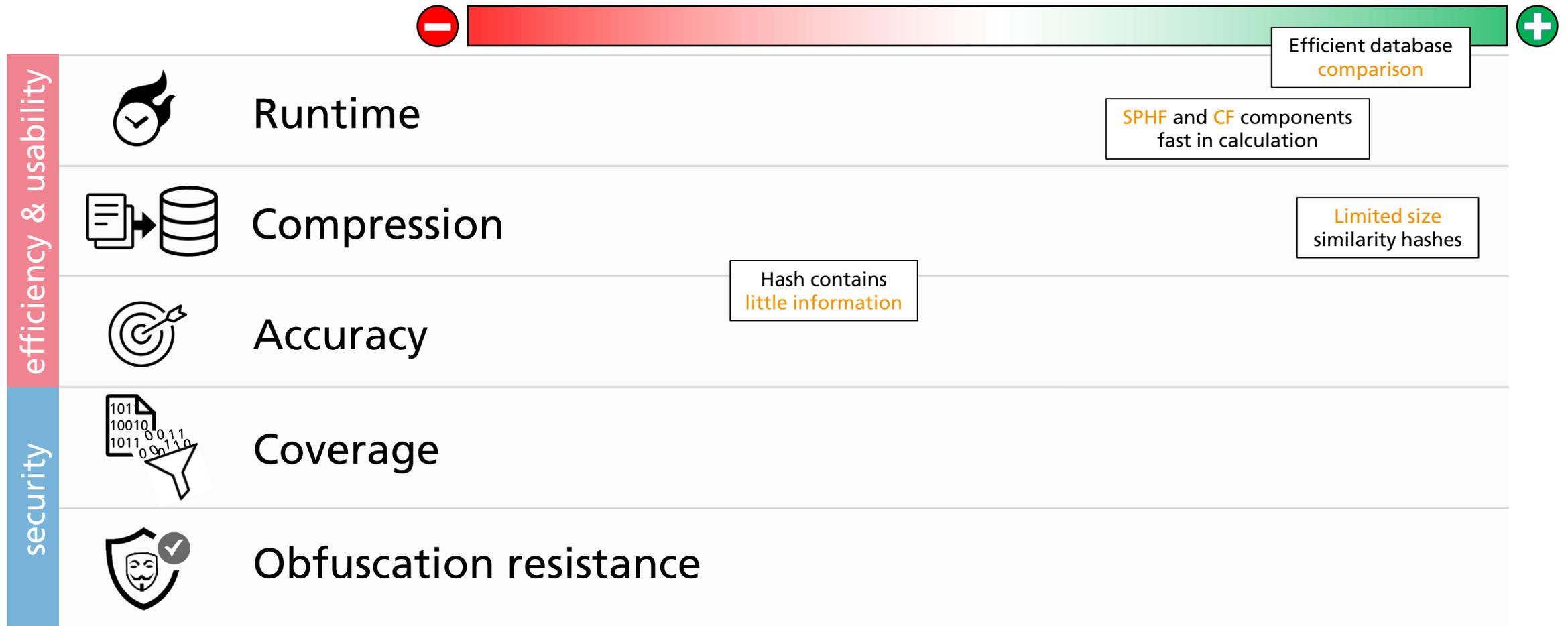
Similarity Hashing – ssdeep – Properties



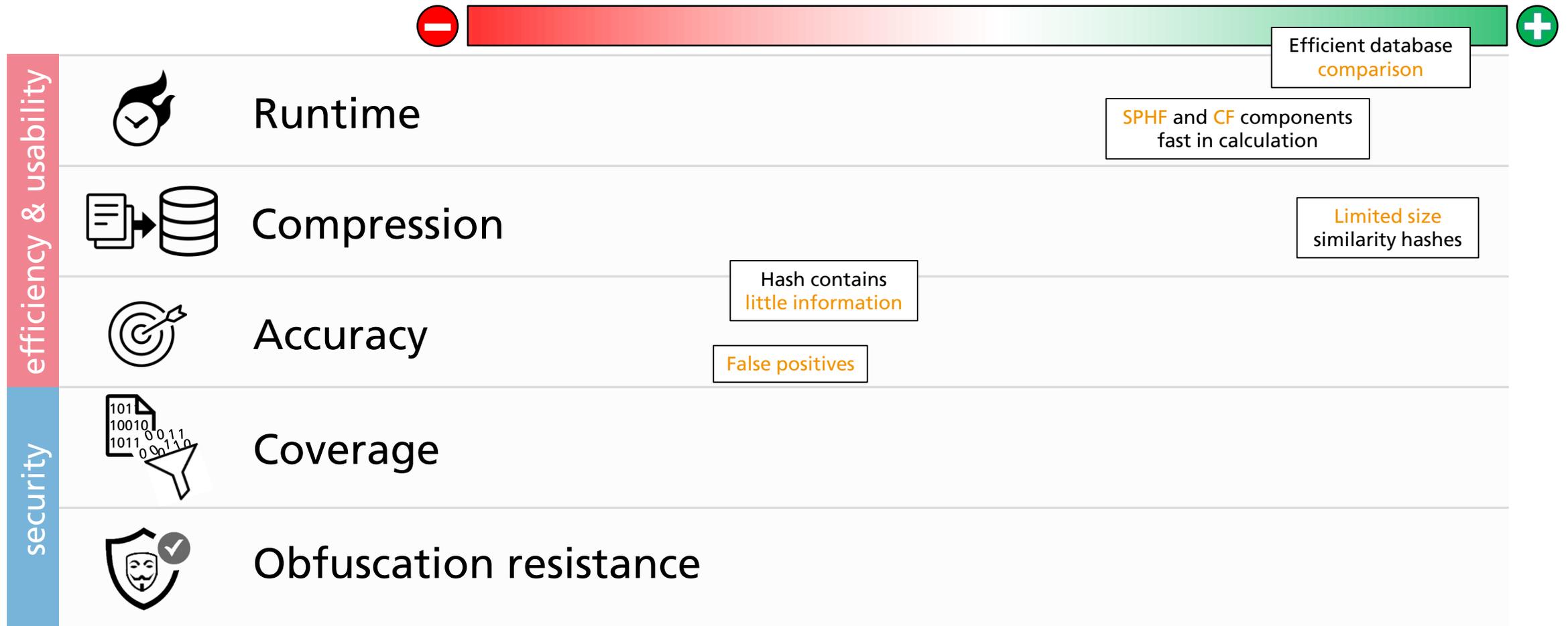
Similarity Hashing – ssdeep – Properties



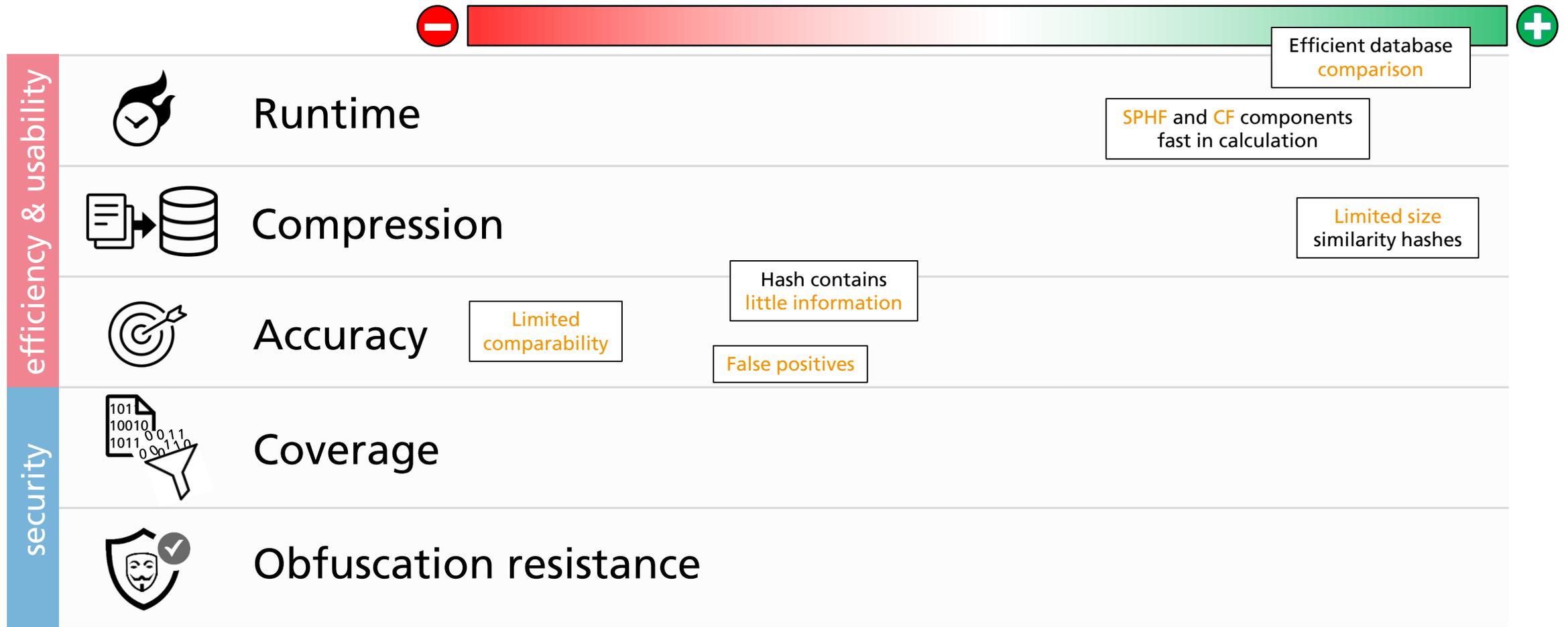
Similarity Hashing – ssdeep – Properties



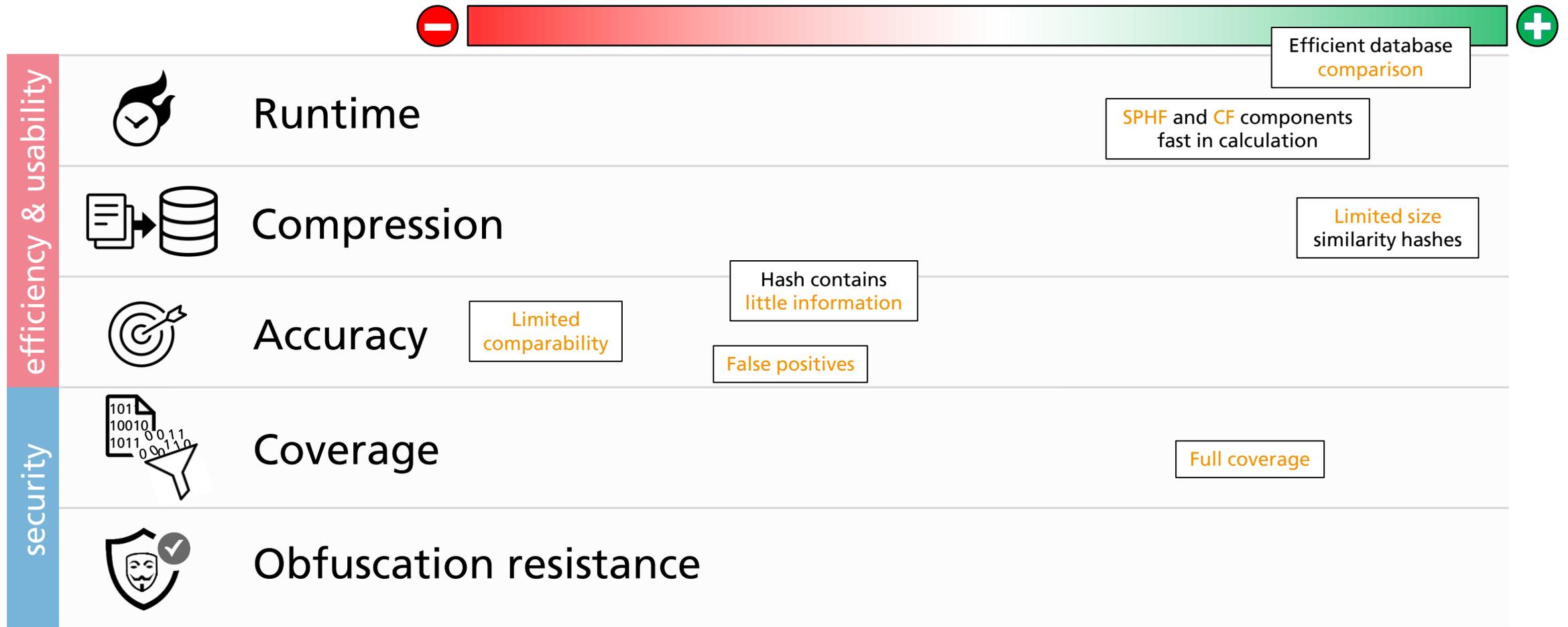
Similarity Hashing – ssdeep – Properties



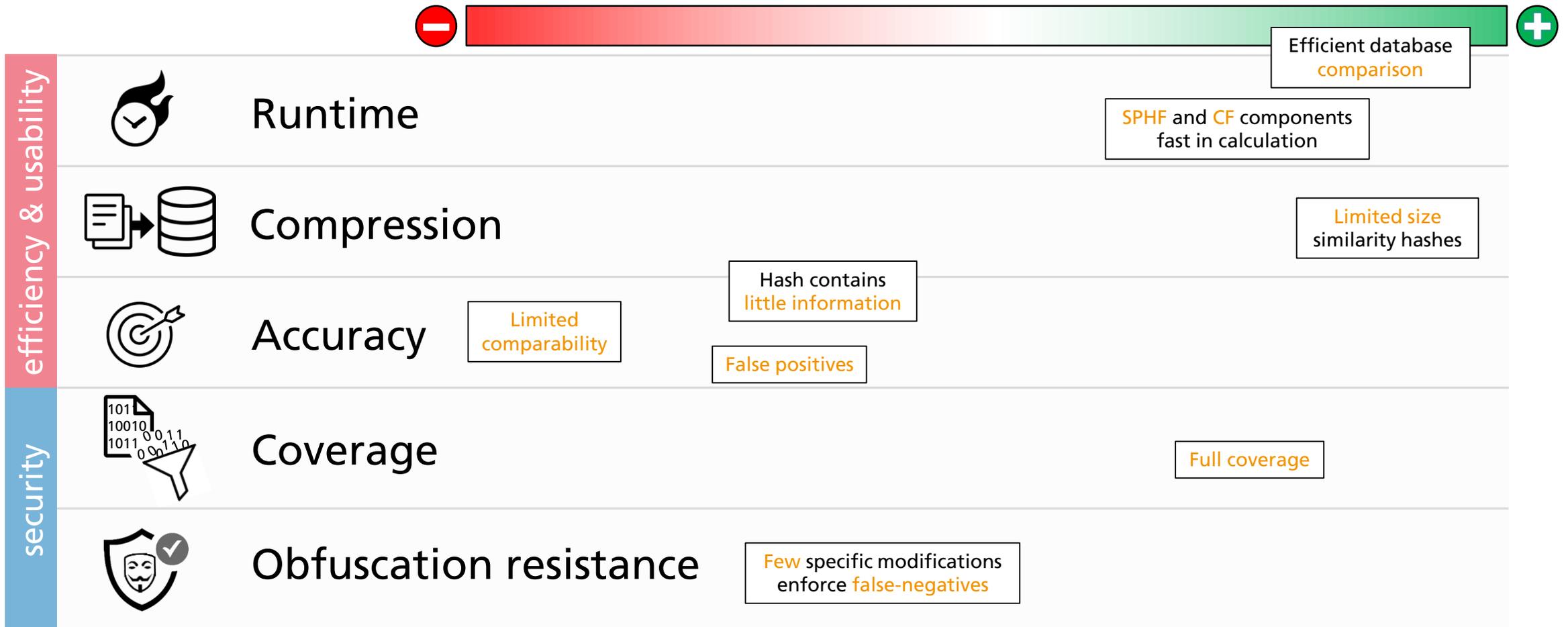
Similarity Hashing – ssdeep – Properties



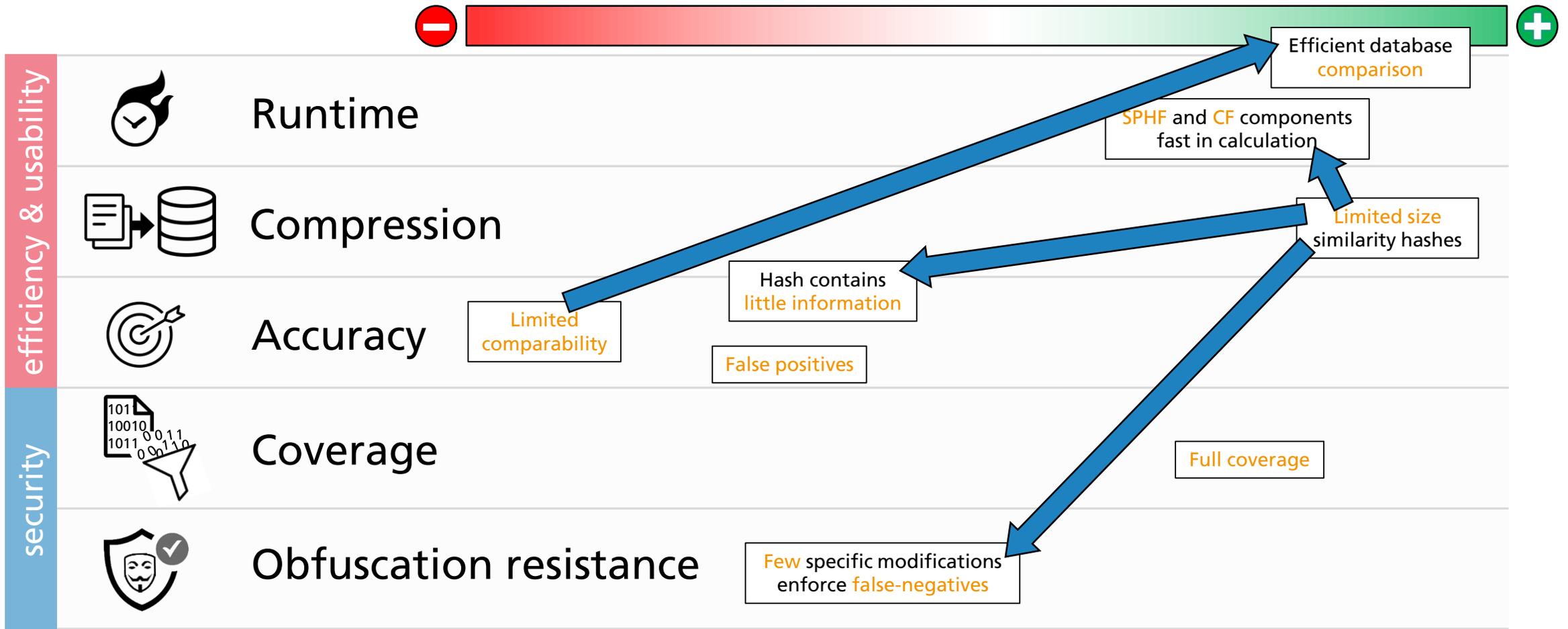
Similarity Hashing – ssdeep – Properties



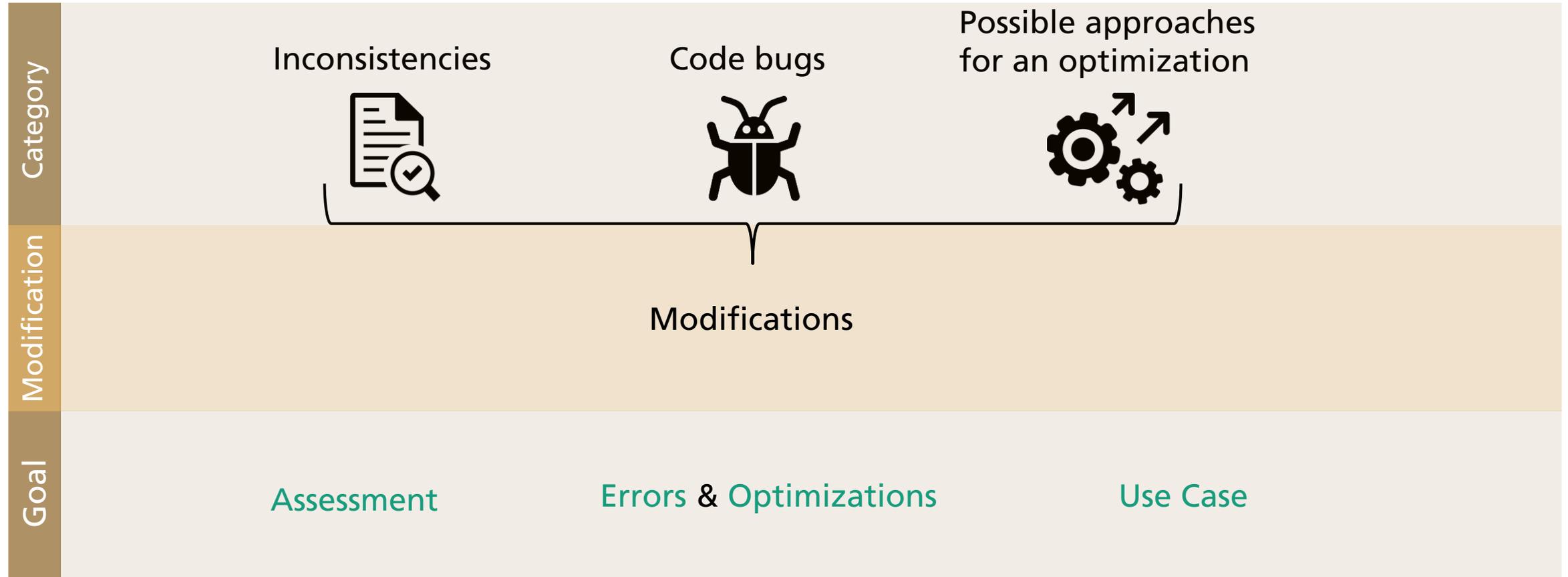
Similarity Hashing – ssdeep – Properties



Similarity Hashing – ssdeep – Properties



Analysis and improvement of ssdeep



Analysis – Inconsistencies

■ Restructuring of the SPHF

- Compatibility: 32 character-limitation

- Modification `-no32lim`

■ CF

- Maximum matching score

- Modification `-nomax`

- Distance calculation via position array

- Modification `-nopa`

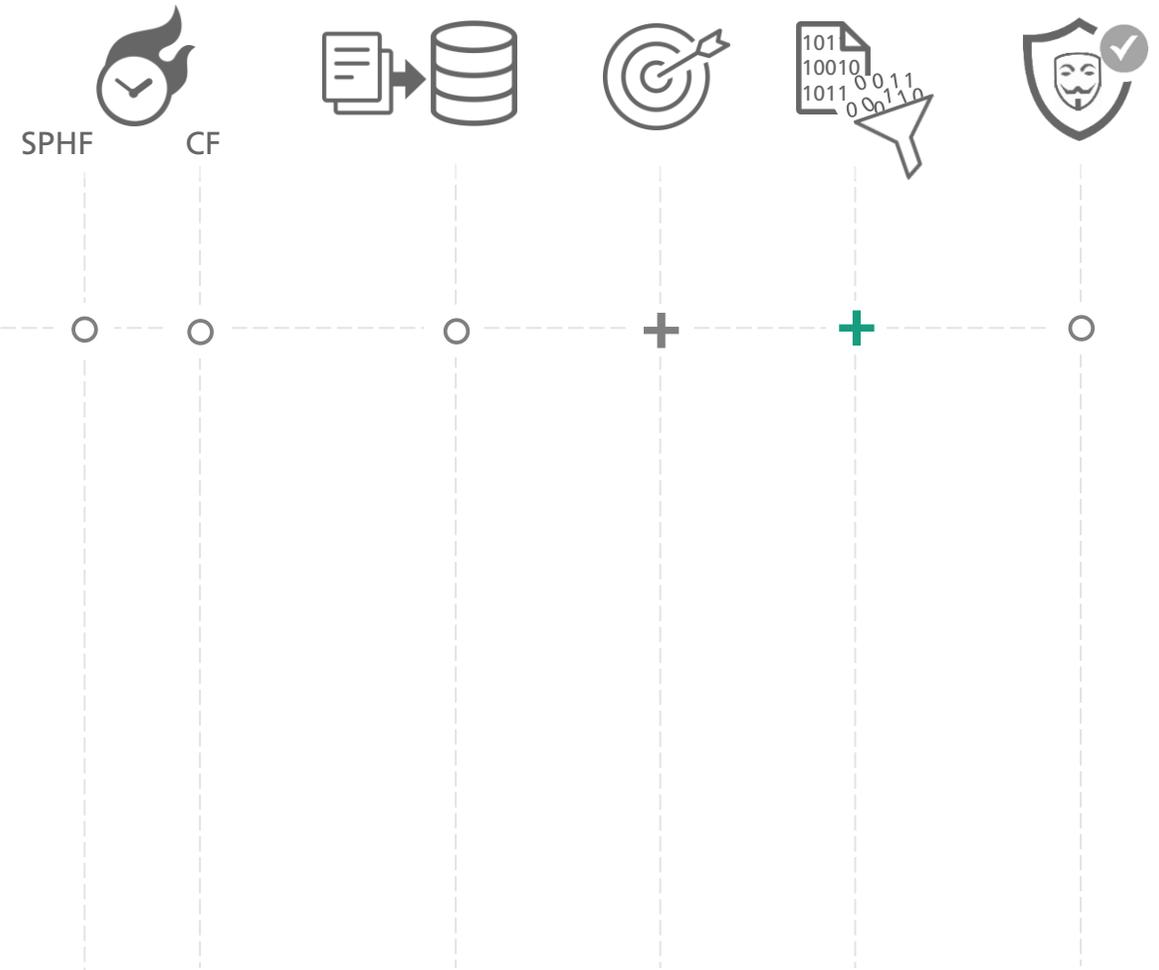
- Common substring

- Modification `-nocommonsub`



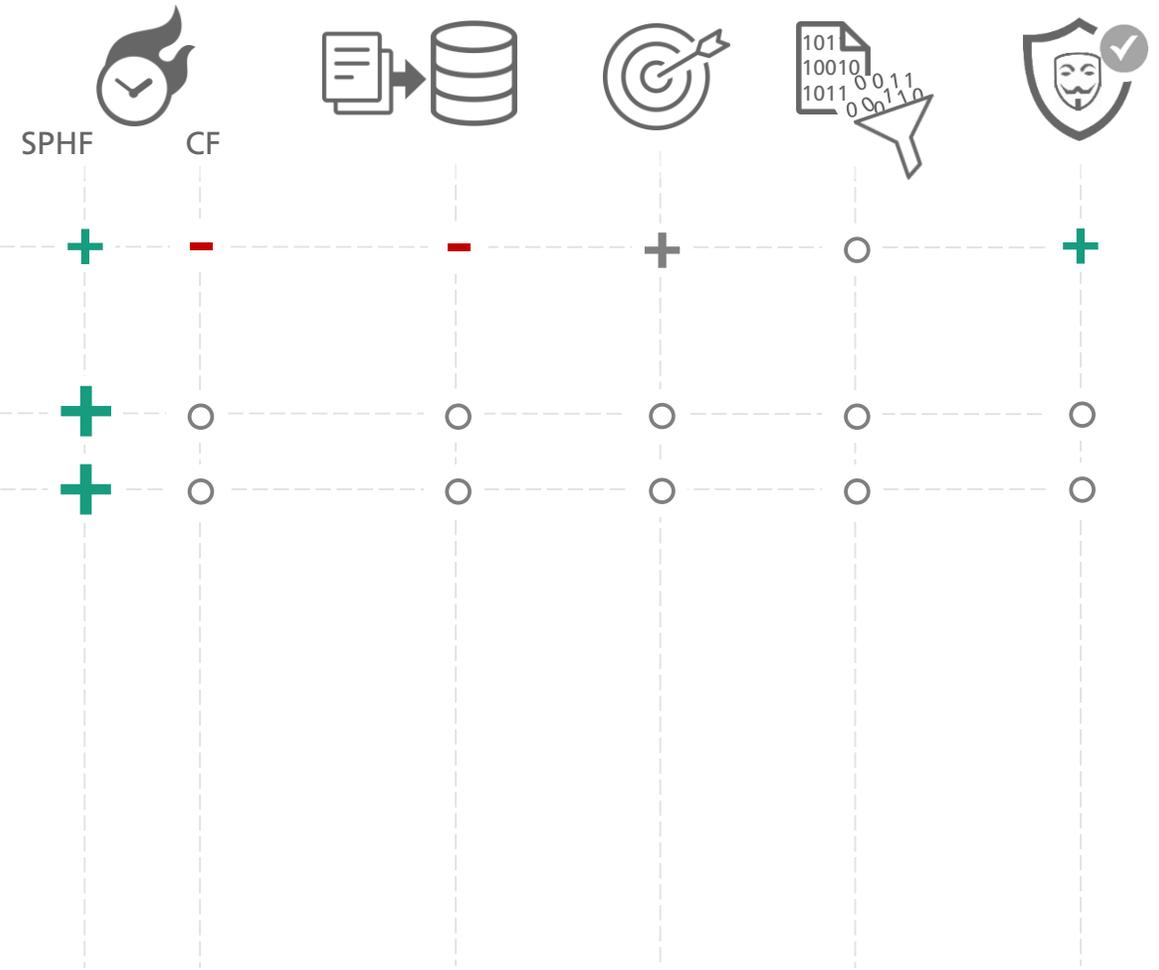
Analysis – Errors

- Last segment bug
 - Empty or missing segment
 - Modification *-bugfix*



Analysis – Optimizations

- Increased trigger progression constant
 - Modification `-4b`
- Exchanging the rolling hash function
 - Modification `-djb2`
 - Modification `-polynomial`



Analysis – Summarized

	SPHF	CF					
 <ul style="list-style-type: none"> Modification <code>-nocommonsub</code> Modification <code>-nomax</code> Modification <code>-nopa</code> 	○	○	○	+	○	+	+
 <ul style="list-style-type: none"> Modification <code>-no32lim</code> Modification <code>-bugfix</code> <p>-refactored</p>	+	-	-	+	+	+	+
 <ul style="list-style-type: none"> Modification <code>-4b</code> Modification <code>-djb2</code> Modification <code>-polynomial</code> 	+	-	-	+	○	+	+
	+	○	○	○	○	○	○

Evaluation

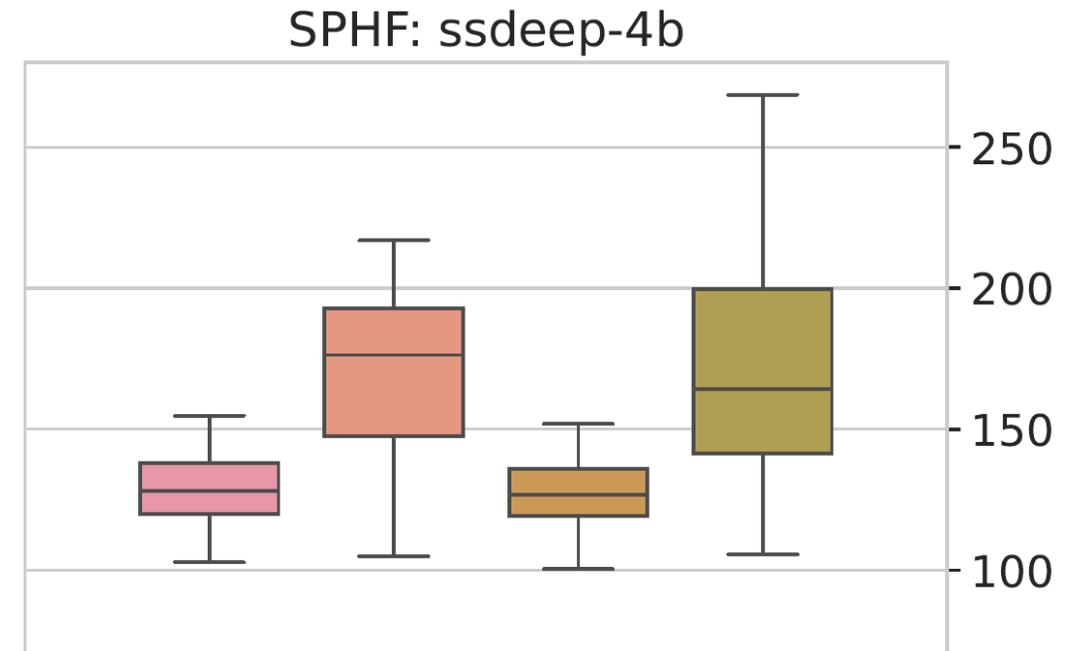
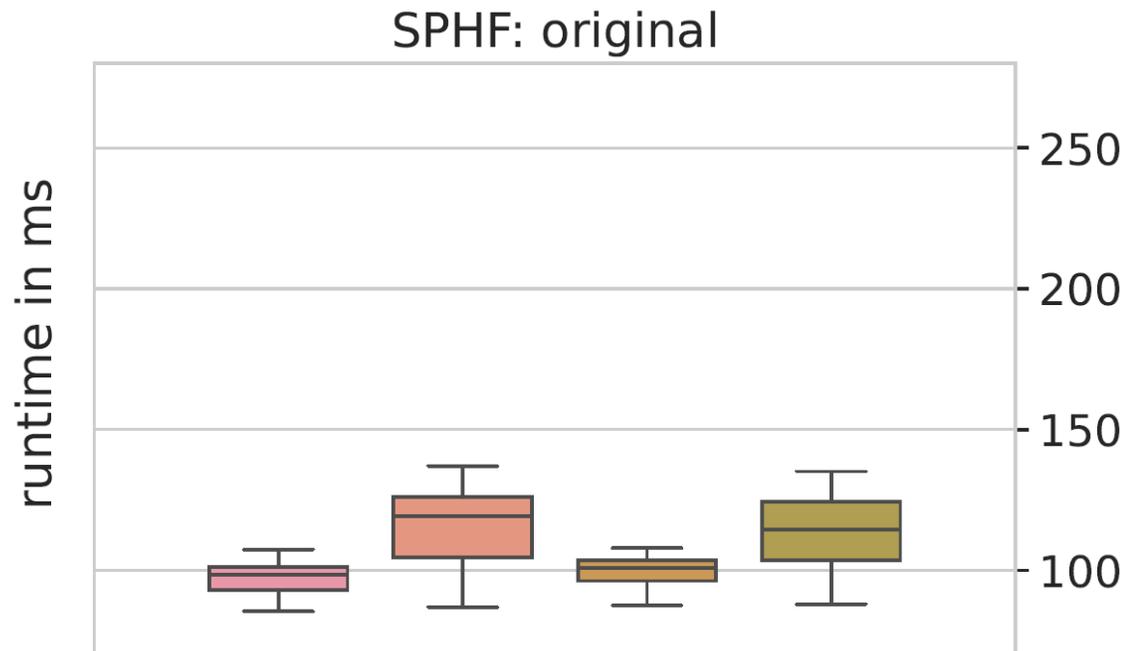
- Similarity Hashing properties
- Most promising candidates
 - SPHF
 - -refactored-djb2
 - -refactored-4b-djb2
 - SPHF & CF
 - -refactored-djb2-nomax
 - -refactored-4b-djb2-nocommonsub
- Publicly available on our GitHub repository:
<https://github.com/fkie-cad/ssdeeper>



Evaluation – Runtime – SPHF

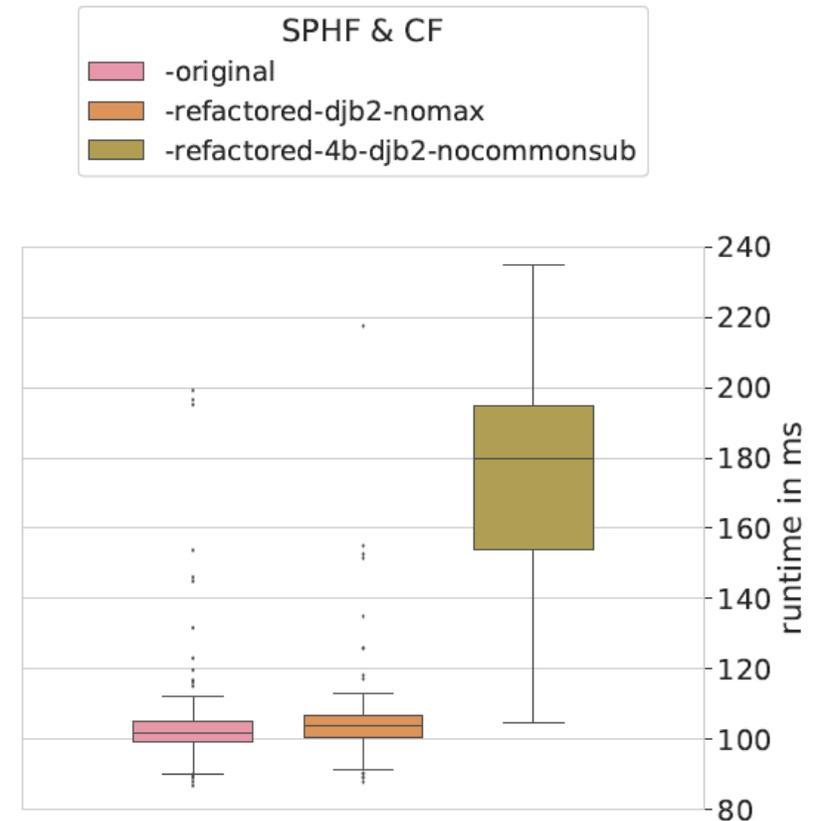
SPHF	Properties of the hash database					Average throughput (KB/ms)	
	Compression ratio	Avg(len(s1))	Avg(len(s2))	Signature ratio	Backward-compatibility	Real data	Synthetic data
-original						48.629	69.619
-refactored-djb2						57.608	89.068
-refactored-4b-djb2						58.876	89.751

Evaluation – Runtime – CF



Evaluation – Runtime – SPHF & CF

SPHF & CF	Average runtime
-original	105.158
-refactored-djb2-nomax	104.826
-refactored-4b-djb2-nocommonsub	175.919



Evaluation – Compression

SPHF	Properties of the hash database					Average throughput (KB/ms)	
	Compression ratio	Avg(len(s1))	Avg(len(s2))	Signature ratio	Backward-compatibility	Real data	Synthetic data
-original	4548.82	50.51	22.65	3.01	n/a	48.629	69.619
-refactored-djb2	4446.78	50.56	25.49	2.33	no	57.608	89.068
-refactored-4b-djb2	3824.28	77.68	19.79	5.23	no	58.876	89.751

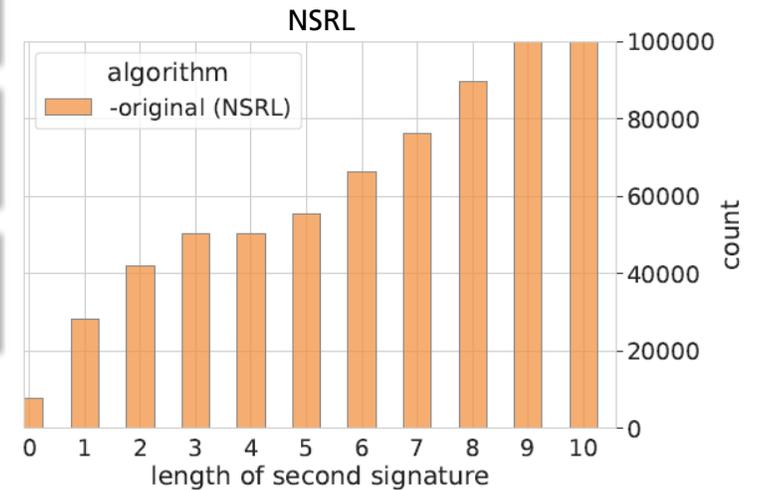
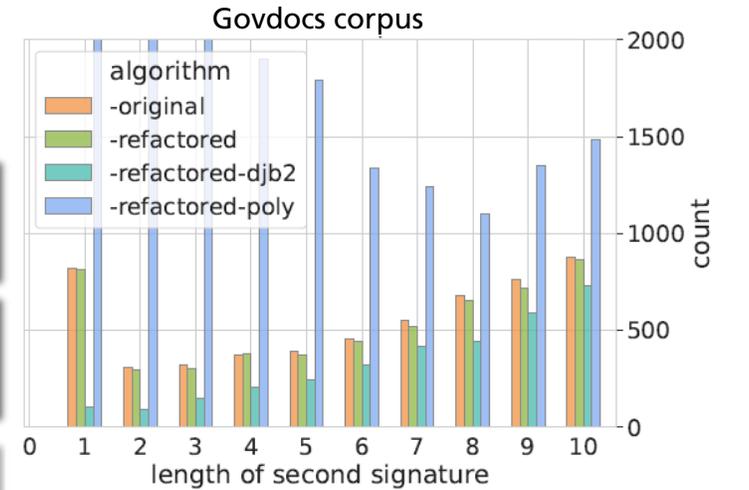
Evaluation – Compression

SPHF	Properties of the hash database					Average throughput (KB/ms)	
	Compression ratio	Avg(len(s1))	Avg(len(s2))	Signature ratio	Backward-compatibility	Real data	Synthetic data
-original	4548.82	50.51	22.65	3.01	n/a	48.629	69.619
-refactored-djb2	4446.78	50.56	25.49	2.33	no	57.608	89.068
-refactored-4b-djb2	3824.28	77.68	19.79	5.23	no	58.876	89.751

Evaluation – Compression

Properties of the hash database

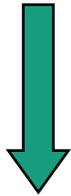
SPHF	Compression ratio	Avg(len(s1))	Avg(len(s2))	Signature ratio	Backward-compatibility
-original	4548.82	50.51	22.65	3.01	n/a
-refactored-djb2	4446.78	50.56	25.49	2.33	no
-refactored-4b-djb2	3824.28	77.68	19.79	5.23	no



Evaluation – Accuracy



synthetic data
corpus

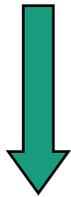


- Randomness provided by input
 - Choice of rolling hash function
- Assessment of modification -4b

Evaluation – Accuracy



synthetic data
corpus



- Randomness provided by input
 - Choice of rolling hash function
- Assessment of modification `-4b`



real data
fragment corpus

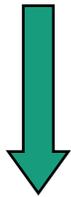


- Worst case scenario
- Assessment of modification `-4b`, `-nocommonsub`, `-djb2` & `-nomax`

Evaluation – Accuracy



synthetic data
corpus



- Randomness provided by input
 - Choice of rolling hash function
- Assessment of modification `-4b`



real data
fragment corpus



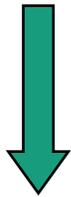
- Worst case scenario
- Assessment of modification `-4b`,
`-nocommonsub`, `-djb2` & `-nomax`

- How to define accuracy and similarity?
 - The smallest error?
 - Here, variance & reliability

Evaluation – Accuracy



synthetic data
corpus



- Randomness provided by input
 - Choice of rolling hash function
- Assessment of modification `-4b`



real data
fragment corpus



- Worst case scenario
- Assessment of modification `-4b`,
`-nocommonsub`, `-djb2` & `-nomax`



real data
step-by-step
corpus

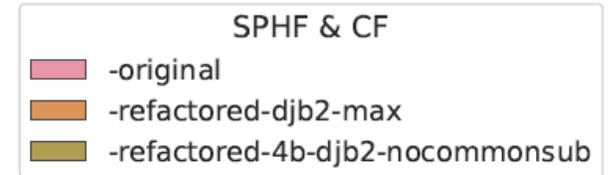
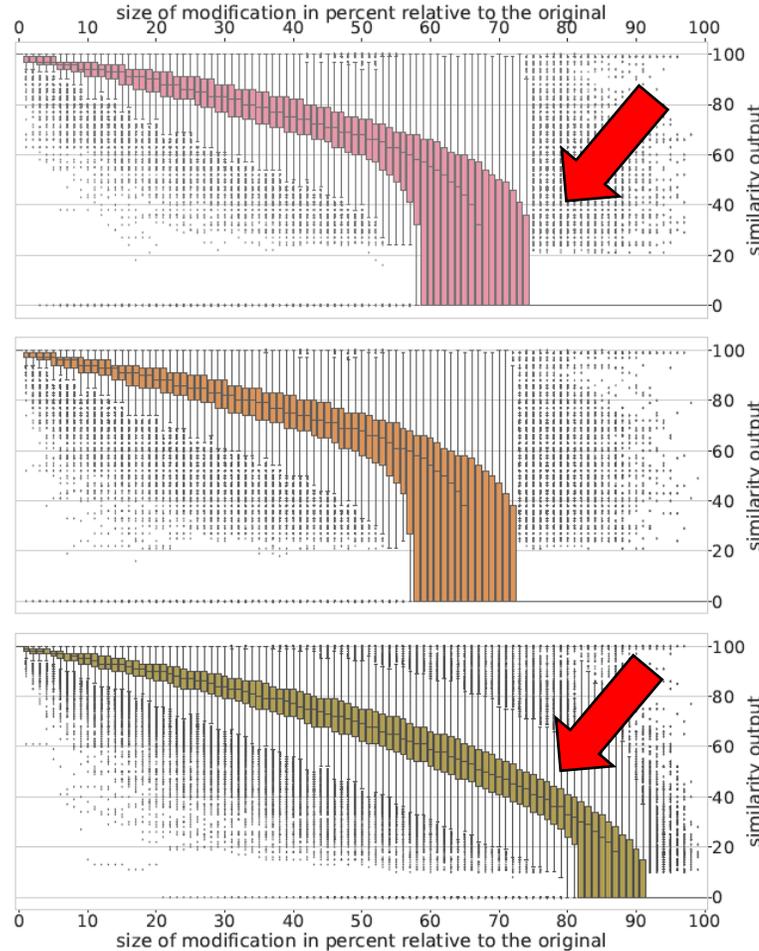
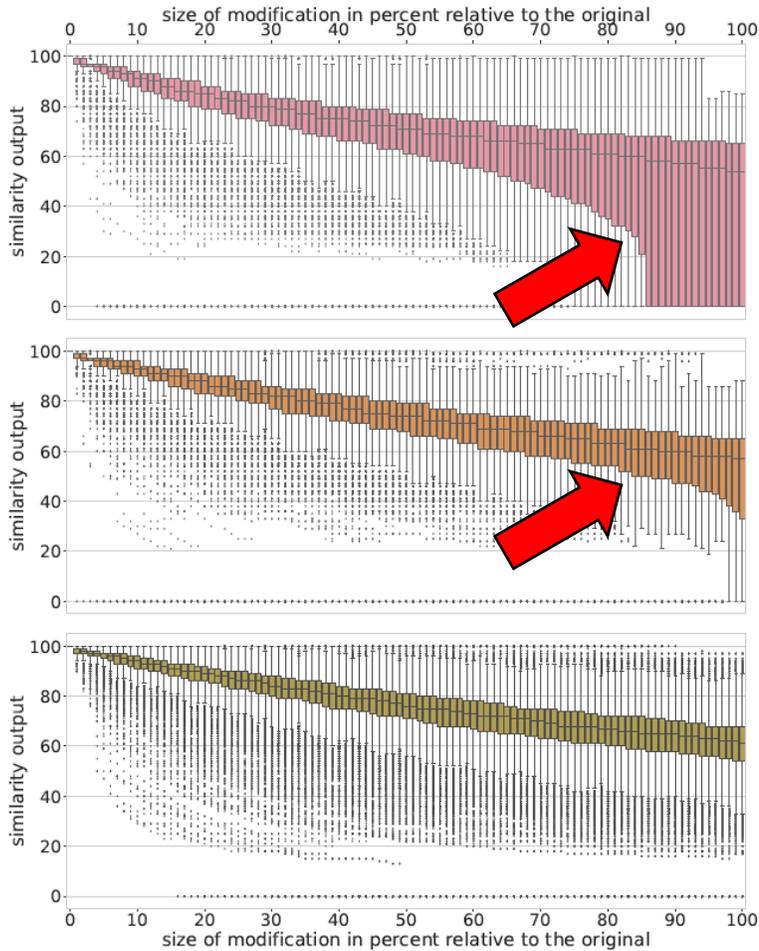


- How to define accuracy and similarity?
 - The smallest error?
 - Here, variance & reliability

Evaluation – Accuracy – Real data

Insertion of bytes

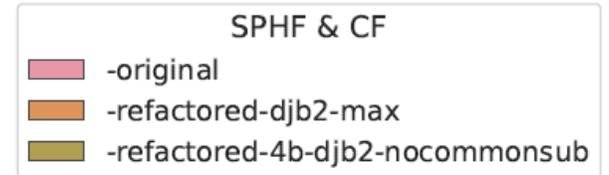
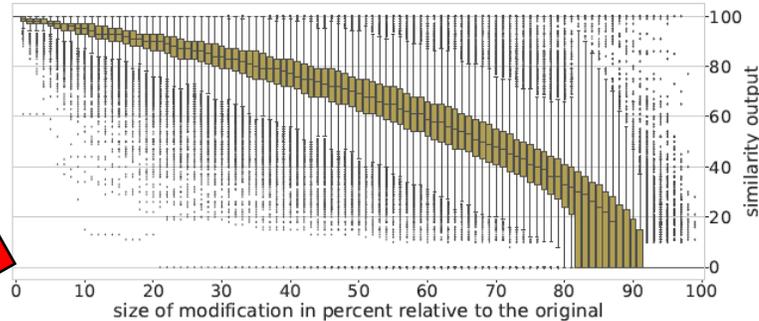
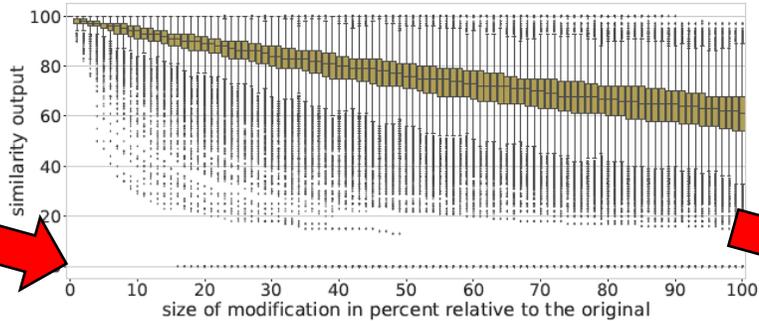
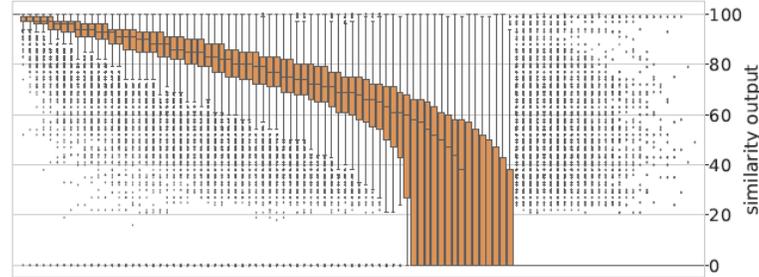
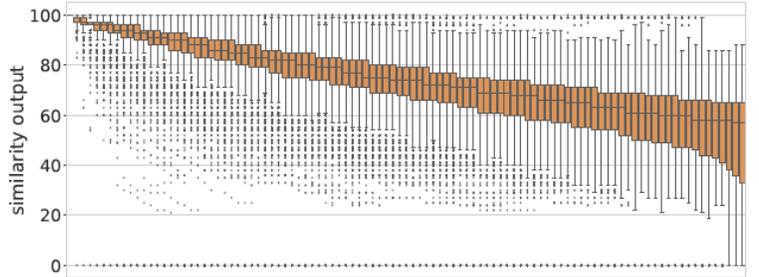
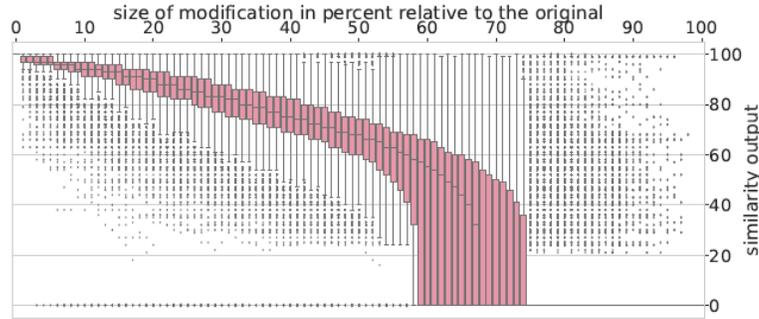
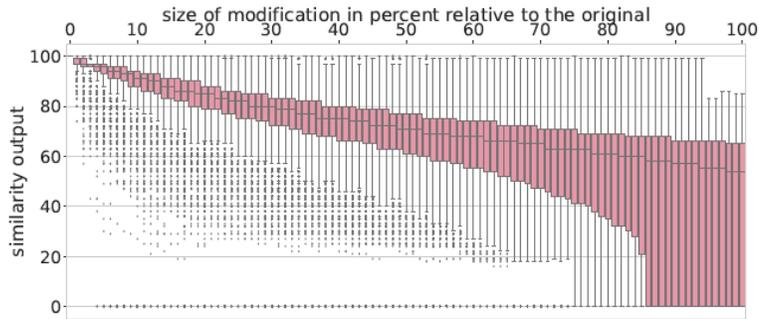
Deletion of bytes



Evaluation – Accuracy – Real data

Insertion of bytes

Deletion of bytes



Evaluation – Security

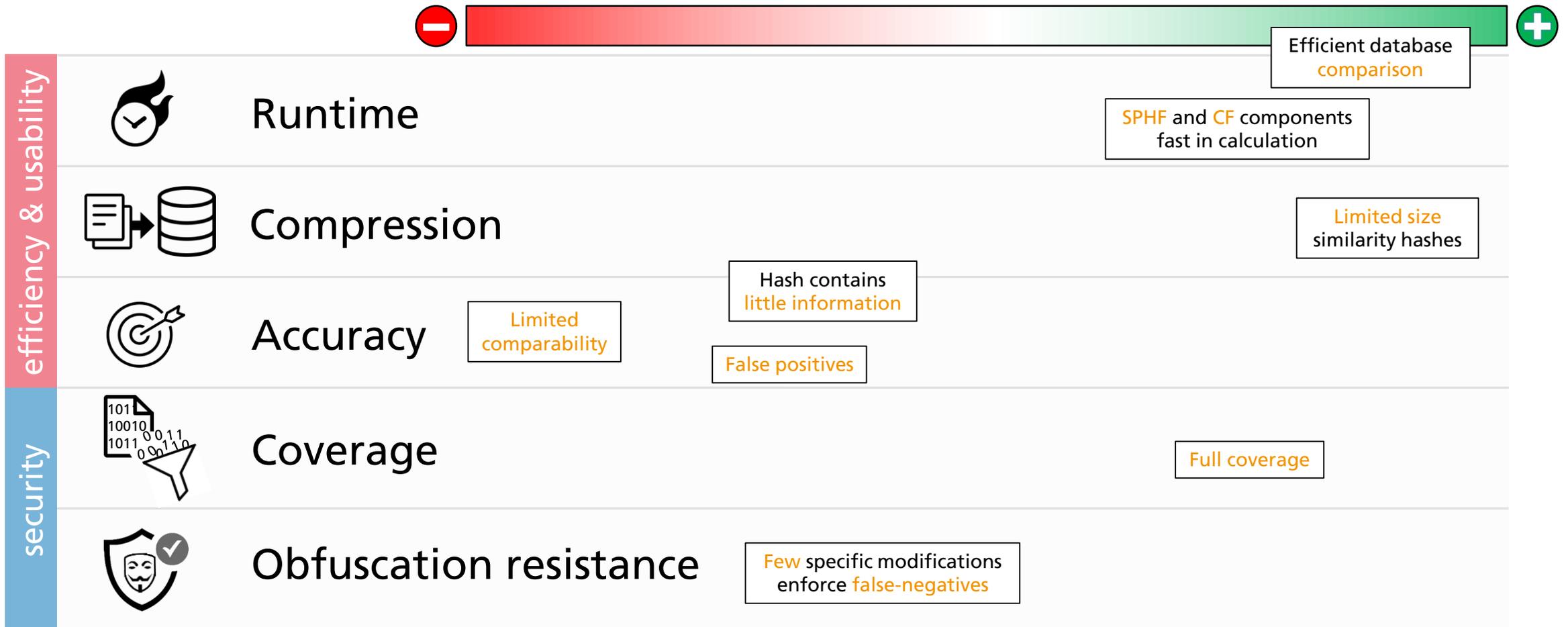
- Obfuscation resistance
 - Possibilities
 - Insertion of trigger points
 - Removal of common substrings
 - ...
 - `-refactored-4b-djb2-nocommonsub`
 - Longer hash values, **yet not really secure**

Evaluation – Security

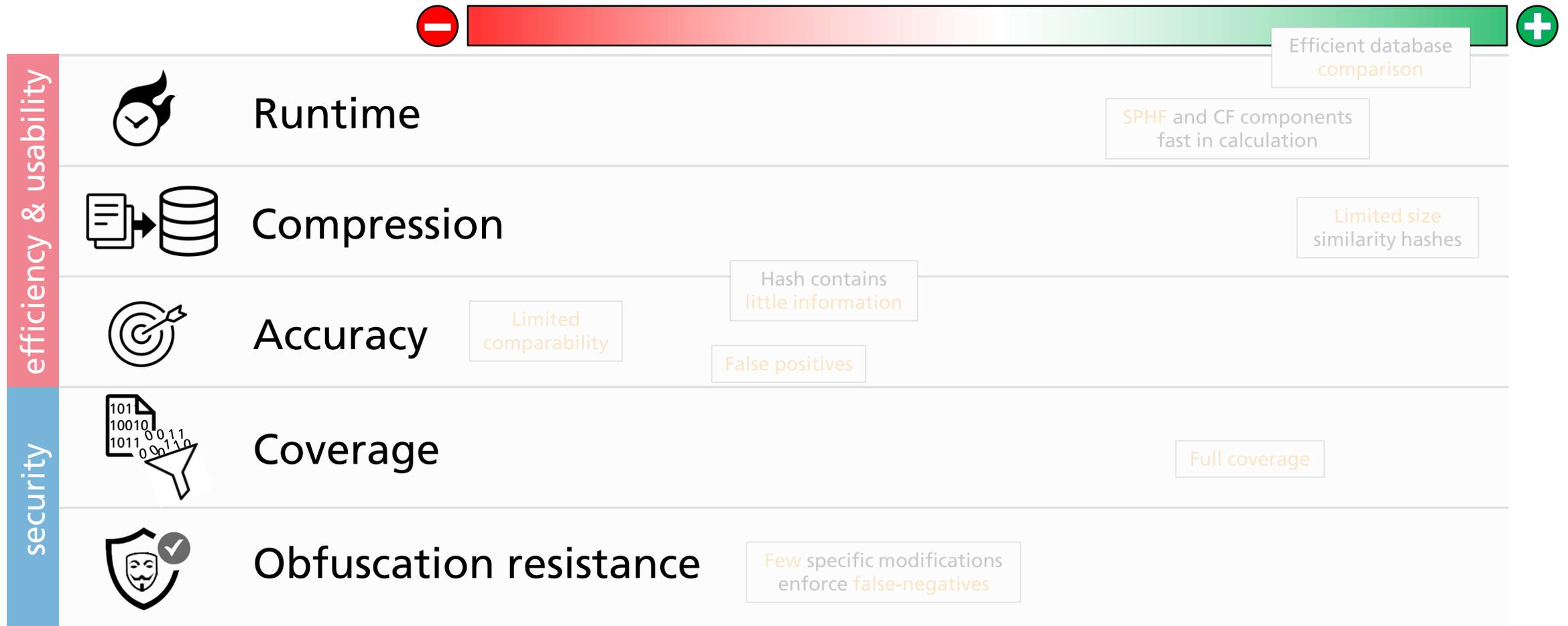
- Obfuscation resistance
 - Possibilities
 - Insertion of trigger points
 - Removal of common substrings
 - ...
 - `-refactored-4b-djb2-nocommonsub`
 - Longer hash values, **yet not really secure**

- Coverage
 - `-bugfix` & `-refactored`

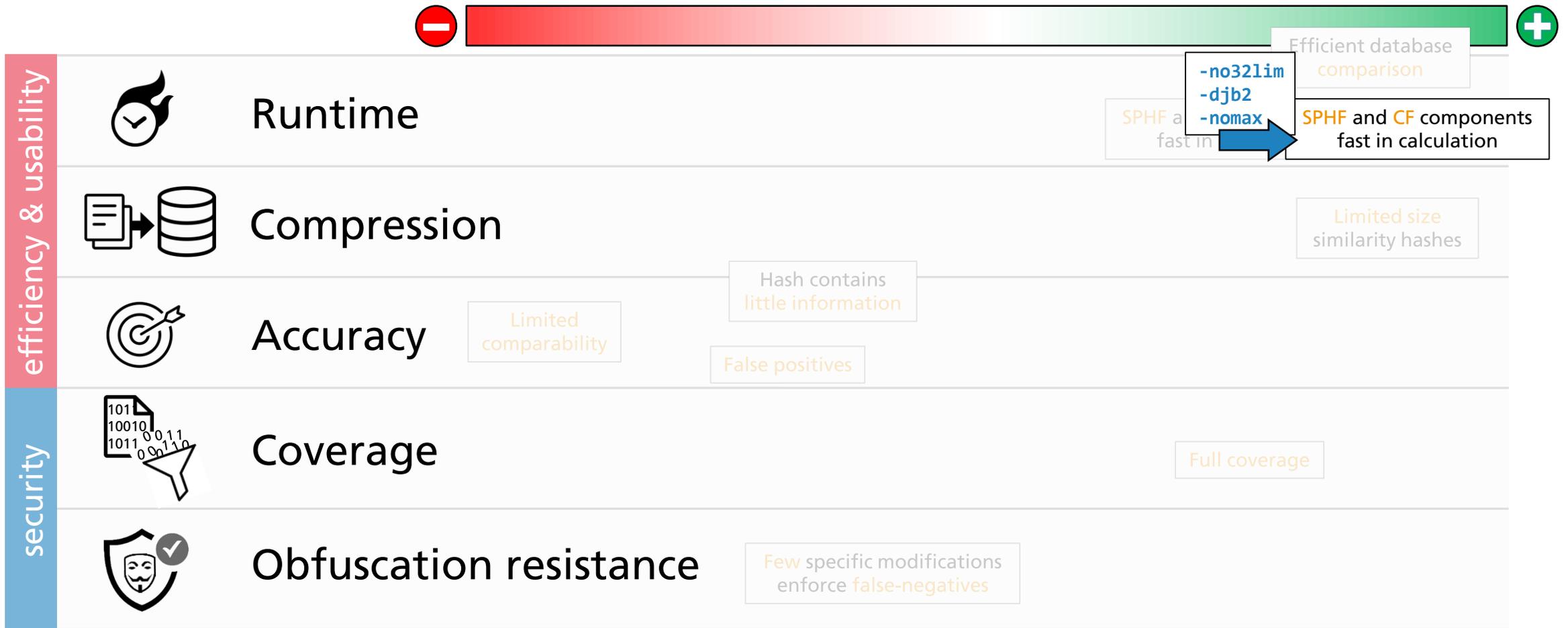
Evaluation – Summary – **ssdeep**



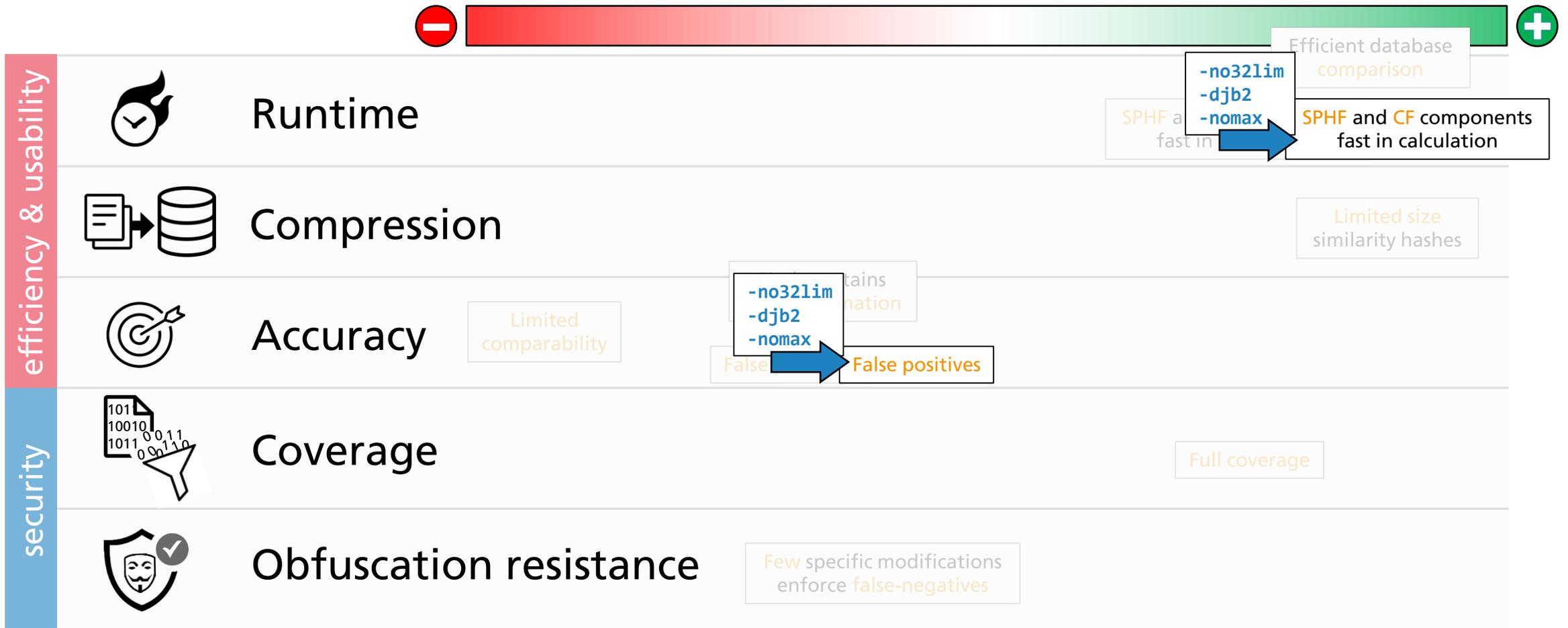
Evaluation – Summary – ssdeep-refactored-djb2-nomax



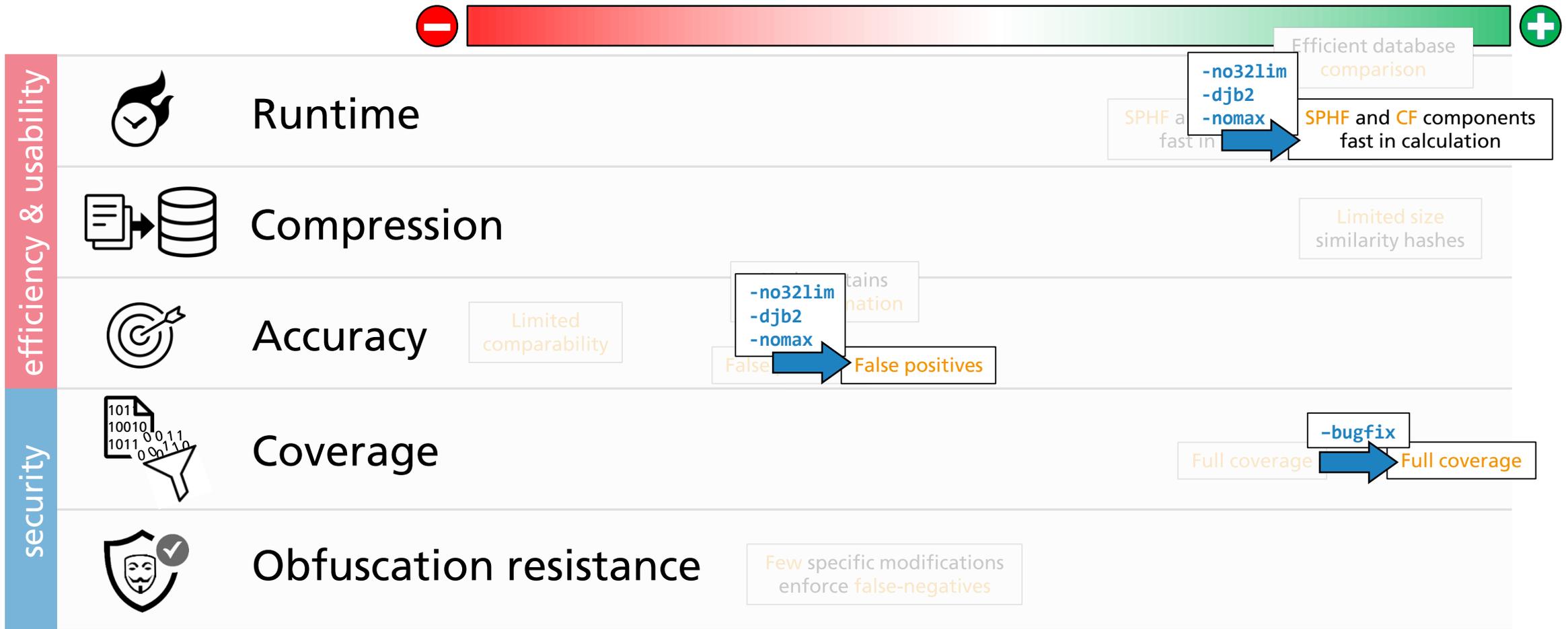
Evaluation – Summary – ssdeep-refactored-djb2-nomax



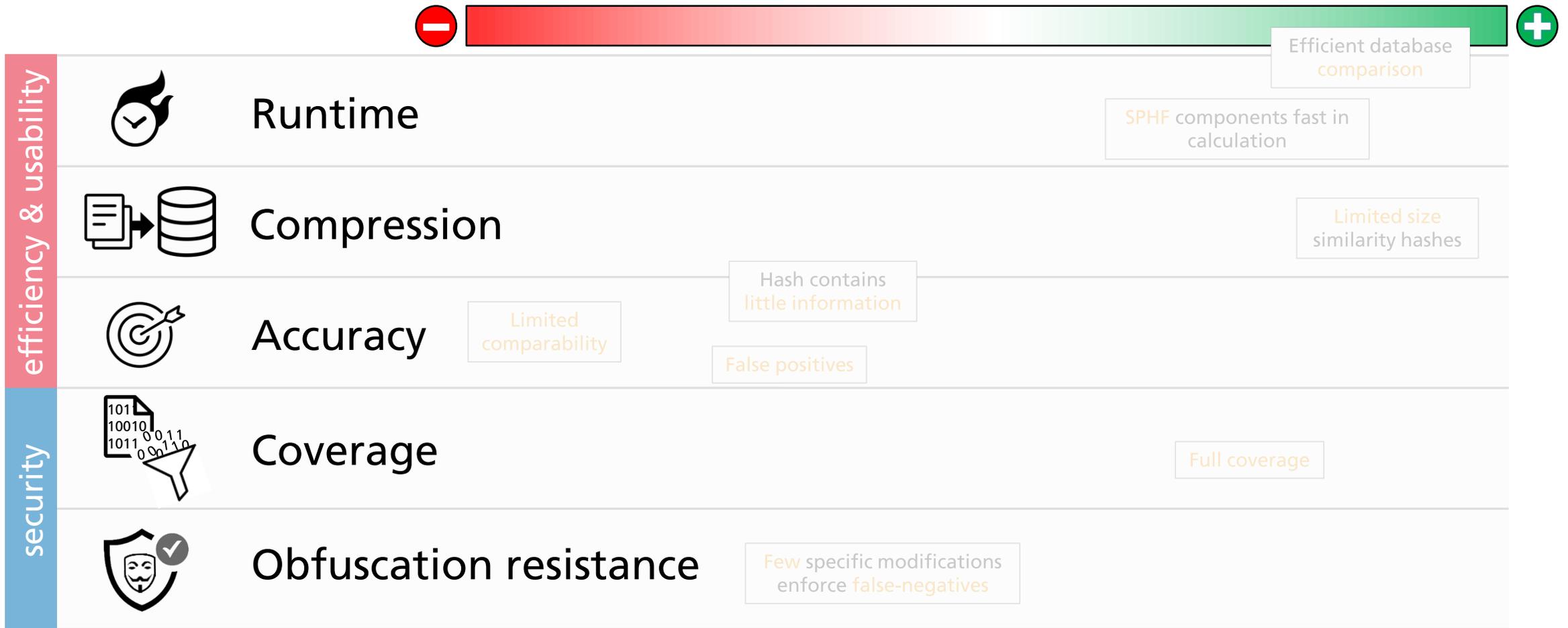
Evaluation – Summary – ssdeep-refactored-djb2-nomax



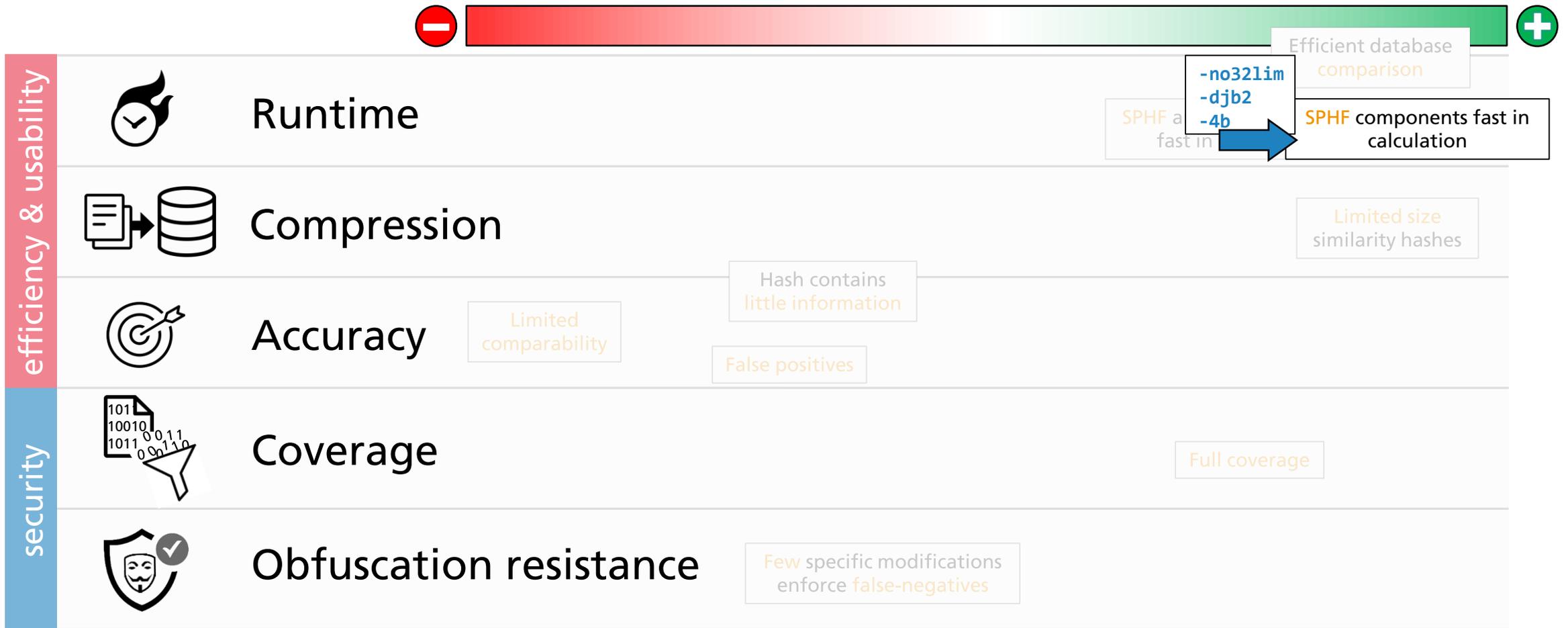
Evaluation – Summary – ssdeep-refactored-djb2-nomax



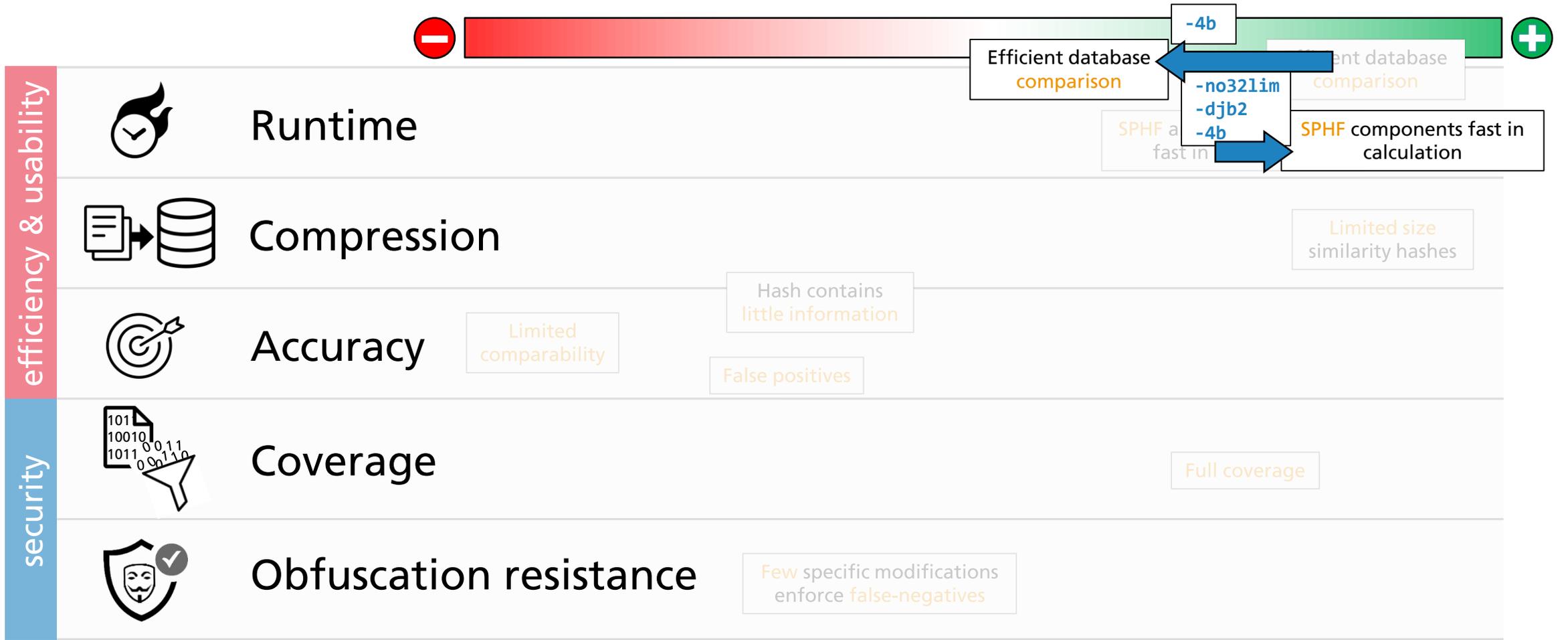
Evaluation – Summary – ssdeep-refactored-4b-djb2-nocommonsub



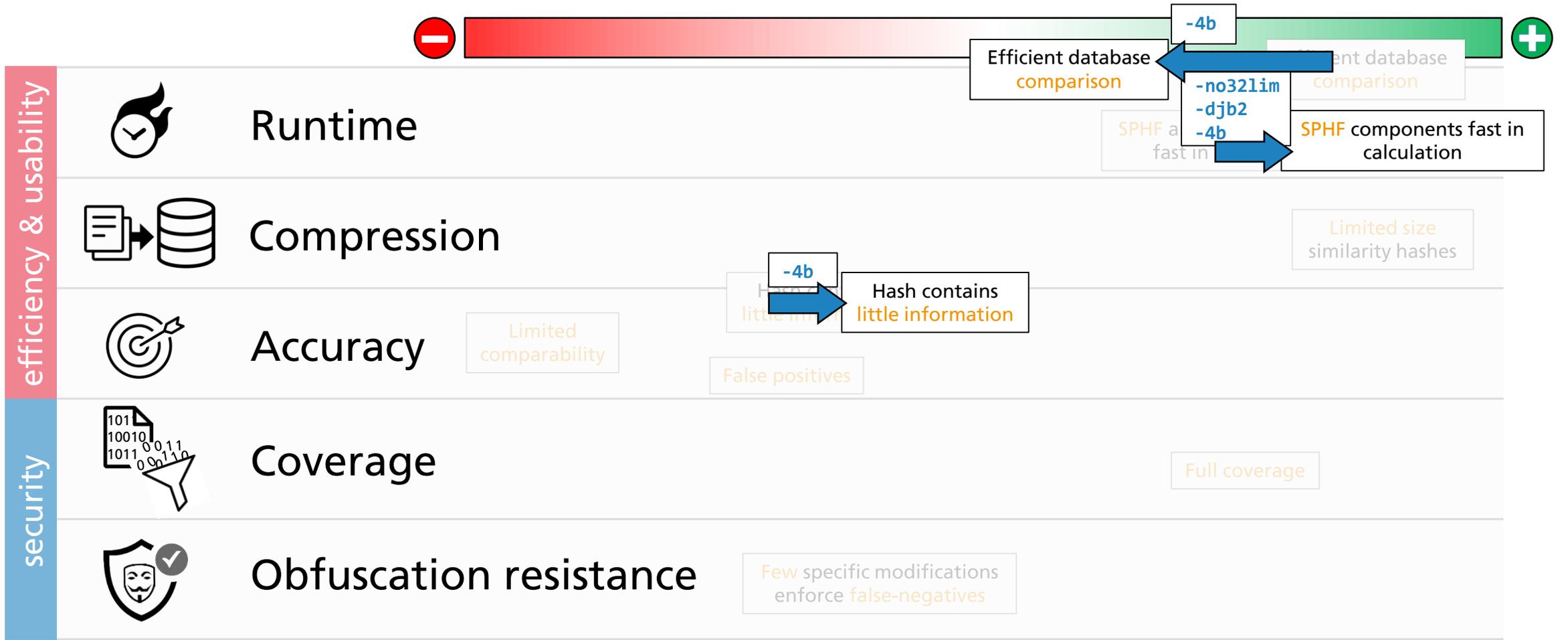
Evaluation – Summary – ssdeep-refactored-4b-djb2-nocommonsub



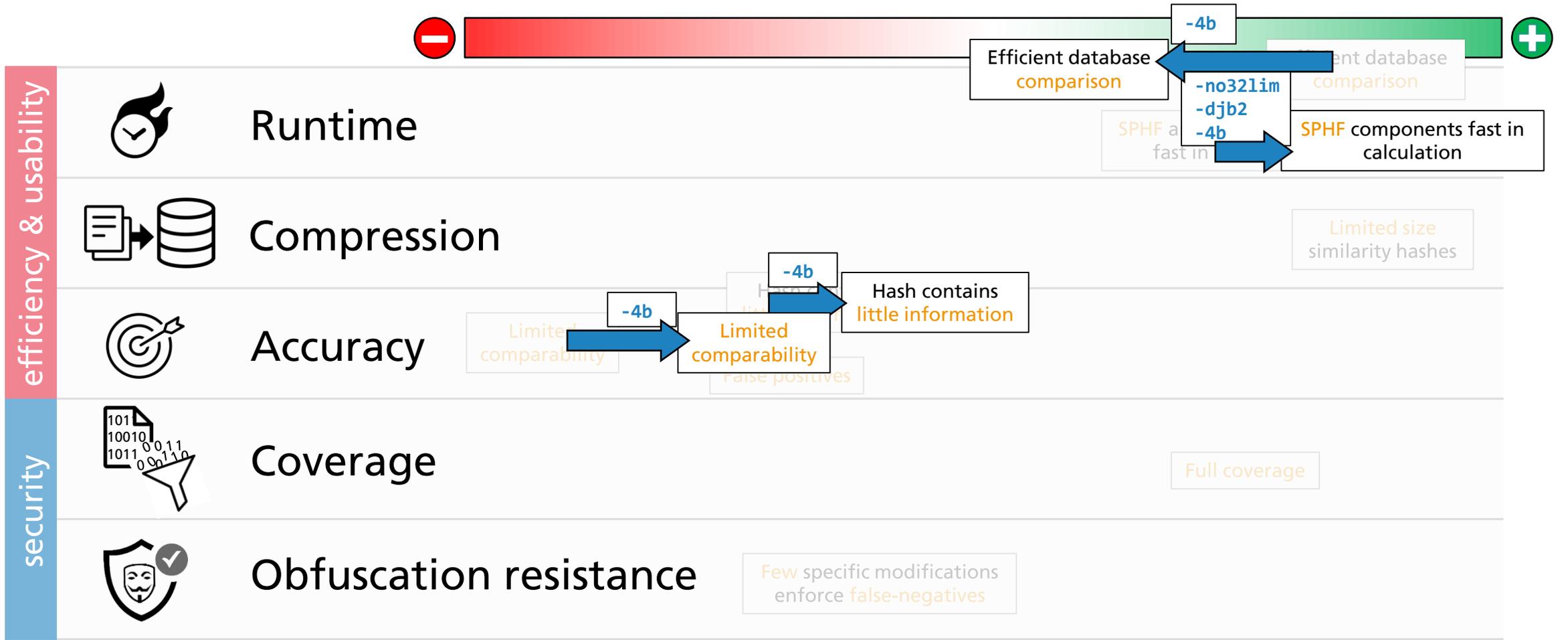
Evaluation – Summary – ssdeep-refactored-4b-djb2-nocommonsub



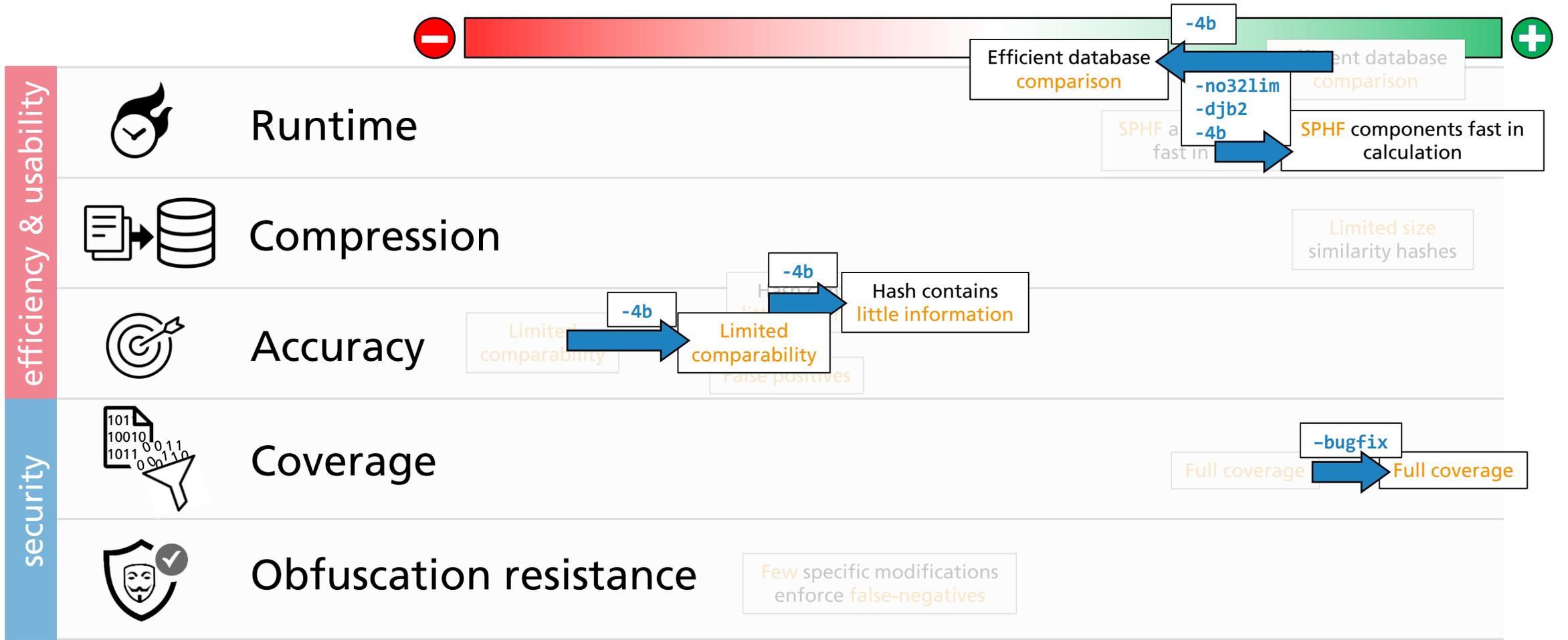
Evaluation – Summary – ssdeep-refactored-4b-djb2-nocommonsub



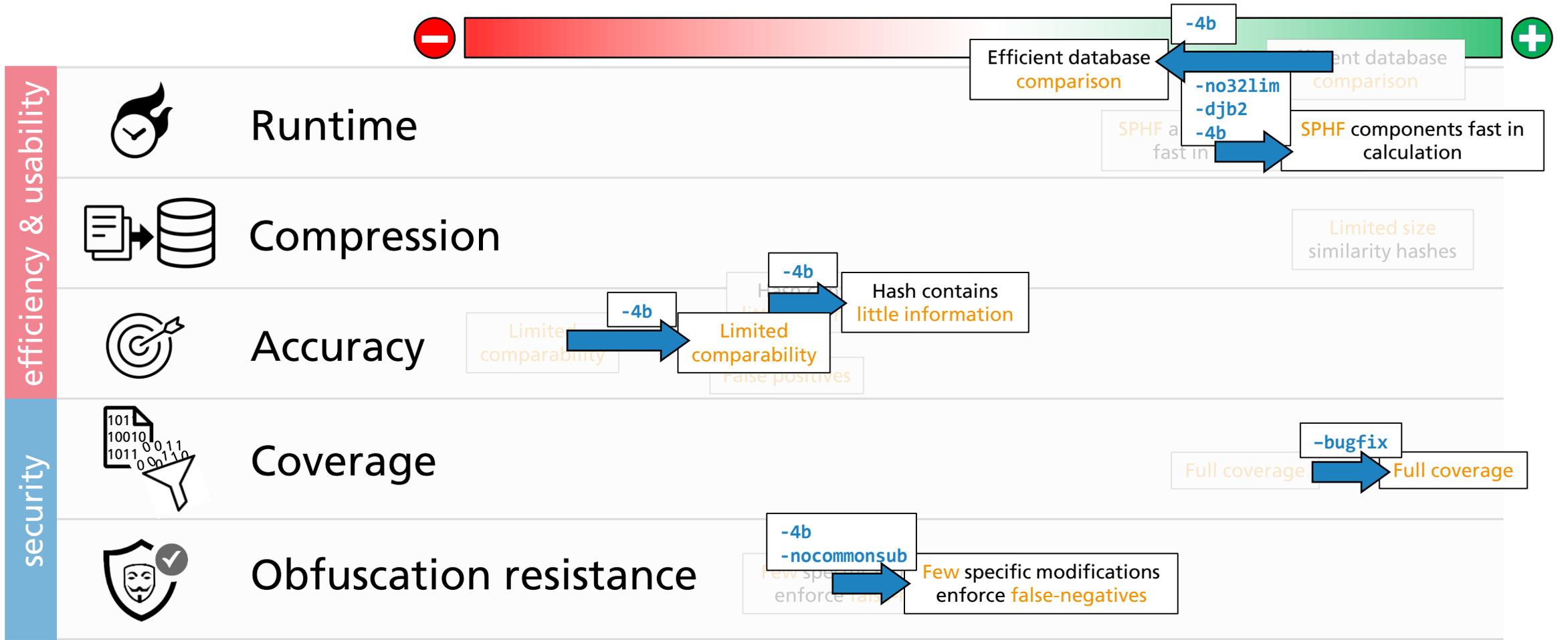
Evaluation – Summary – ssdeep-refactored-4b-djb2-nocommonsub



Evaluation – Summary – ssdeep-refactored-4b-djb2-nocommonsub



Evaluation – Summary – ssdeep-refactored-4b-djb2-nocommonsub



Conclusion

- Assessment and Evaluation of
 - ssdeep
 - proposed modifications
- Publicly available:
 - Modifications and Implementations
 - Data corpus
 - Evaluation results and scripts
 - Files highlighting issues

- `ssdeep-refactored-djb2-nomax`
 - More efficient SPHF and CF
 - Fewer false positives
 - Use case: **Filtering** and **minimizing** extraneous data

- `ssdeep-refactored-4b-djb2-nocommonsub`
 - Slower but larger range of comparisons
 - More efficient SPHF
 - Fewer false negatives
 - Use case: **Detection** of **modified suspicious files**

Thank you!
Any Questions?

Contact Information

Carlo Jakobs
carlo.jakobs@fkie.fraunhofer.de

Martin Lambertz
martin.lambertz@fkie.fraunhofer.de

Jan-Niclas Hilgert
jan-niclas.hilgert@fkie.fraunhofer.de

Code and data

<https://github.com/fkie-cad/ssdeeper>