# A Data Mining Perspective on Explainable AIOps with Applications to Software Maintenance

Presented by **Youcef REMIL**

Proposed by **Infologic R&D**

Advised by **Pr. Jean-François BOULICAUT**
**Dr. Mehdi KAYTOUE**
**Dr. Anes BENDIMERAD**

**June 04, 2024**

# Introduction and Motivation – CIFRE Thesis

**1982** - Infologic foundation

**2016** – Infologic R&D Initiated

**2019** – Preventive Maintenance Project

**2020** – Ph.D Thesis on AIOps

ERP Software Editor **Copilote**

- Significant annual growth
- More than **600 sites**
- Over **200K workstations**

# Introduction and Motivation – CIFRE Thesis

**2016** – Infologic R&D Initiated

Collect and storage of **telemetry data**

- Boosting **efficiency/reliability**
- Service **quality**
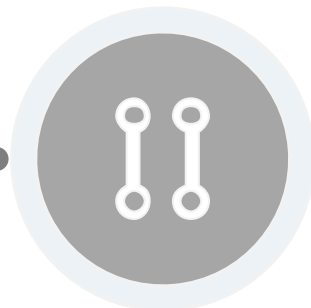- Need for **automation**

# Introduction and Motivation – CIFRE Thesis

**1982** - Infologic foundation

**2016** – Infologic R&D Initiated

**2019** – Preventive Maintenance Project

**2020** – Ph.D Thesis on AIOps

- **Data-centric** approach
- **Real-time** monitoring
- **Proactive** maintenance
- **AIOps\*** and Automation

*__AIOps: AI__ for __Op__erating __S__ystems [Pankaj Prasad and Charley Rich. Market guide for AIOps platforms]

# Introduction and Motivation – CIFRE Thesis

**1982** - Infologic foundation

**2016** – Infologic R&D Initiated

**2019** – Preventive Maintenance Project

**2020** – Ph.D Thesis on AIOps

- Study of **AIOps** field
- **Limitations** of AIOps
- **Development** of effective AIOps solutions
- **Applicative** and **Research Contributions**

*__AIOps: AI__ for **Op**erating **S**ystems [Pankaj Prasad and Charley Rich. Market guide for AIOps platforms]

# Introduction and Motivation

❑ Real pain points of maintenance routines at Infologic

➤ **Lack of standardized and automated maintenance routines with higher costs**

- ▪ Relying mostly on corrective maintenance



- ▪ Example: A detectable **memory leak** at a customer's premises (with **+€450m** annual revenue) **blocked** the departure of all delivery trucks from a factory for **30 minutes**.

# Introduction and Motivation

❑ Real pain points of maintenance routines at Infologic

▶ **Lack of standardized and automated maintenance routines with higher costs**

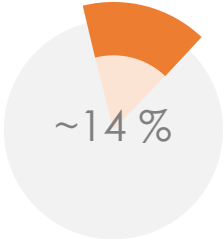- Higher human and resource costs [statistics by the end of 2019]

| Code | Libellé | 01/01/19 00:00 - 31/12/19 00:00 | | | 01/01/20 00:00 - 31/12/20 00:00 | | |
|---|---|---|---|---|---|---|---|
| | | Durée ▼ | Durée moy. | Prod. horaire | Durée | Durée moy. | Prod. hora |
| > 448719 | | 367j 4h 47m | 1h 37m 29s | 0.62 | 361j 5h 19m | 1h 58m 49s | 0. |
| > 604120 | | 166j 5h 18m | 1h 20m 43s | 0.74 | 106j 7h 46m | 1h 19m 51s | 0. |
| > 160249 | | 119j 7h 1m | 1h 34m 1s | 0.64 | 237j 48m 1s | 1h 50m 10s | 0. |
| > 280110 | | 116j 6h 53m | 57m 10s | 1.05 | 98j 1h 9m | 1h 3m 29s | 0. |
| > 091225 | | 108j 7h 43m | 1h 42m 9s | 0.59 | 57j 5h 28m | 1h 27m 20s | 0. |
| > 091790 | | 100j 2h 11m | 1h 35m 52s | 0.63 | 77j 7h 46m 1s | 1h 34m 30s | 0. |
| > 484270 | | 94j 7h 54m | 1h 37s | 0.99 | 81j 4h 56m | 1h 7m 32s | 0. |
| > 091730 | | 93j 7h 22m | 1h 32m 34s | 0.65 | 40j 1h 27m | 1h 35m | 0. |
| > 800000 | | 93j 7h 19m | 56m 33s | 1.06 | 75j 25m | 56m 22s | 1. |
| > 091780 | | 84j 44m | 58m 19s | 1.03 | 78j 18m | 1h 7m 44s | 0. |
| > 840130 | | 80j 6h 48m | 53m 22s | 1.12 | 77j 5h 15m | 56m 13s | 1. |
| > 320557 | | 78j 1h 36m | 1h 18m 51s | 0.76 | 103j 1h 23m 1s | 1h 16m 39s | 0. |
| > 554020 | | 71j 3h 11m | 58m 5s | 1.03 | 31j 7h 51m | 55m 13s | 1. |
| > 724299 | | 68j 3h 41m | 1h 10m 49s | 0.85 | 57j 6h 5m | 1h 12m 57s | 0. |
| > 040340 | | 64j 3h 59m | 1h 21m 54s | 0.73 | 34j 3h 49m | 1h 17m 19s | 0. |
| > 440500 | | 44j 6h 39m | 1h 1m 8s | 0.98 | 44j 3h 49m | 54m 27s | 1. |
| > 200440 | | 44j 5h 39m | 1h 11m 46s | 0.84 | 41j 5h 39m | 1h 18m 11s | 0. |
| > 247301 | | 44j 2h 51m | 1h 42m 51s | 0.58 | 32j 3h 38m | 1h 16m 21s | 0. |
| Total | | 5592j 5h 41m | 1h 2m 48s | 0.9 | 5739j 5h 56m 9s | 1h 7m 44s | 0. |

~5600 days — Maintenance time

~28 pers — Full-time employees

~14 % — Workforce Percentage

# Introduction and Motivation

❑ Real pain points of maintenance routines at Infologic

➤ **Inefficient incident triage and classification**

- Need for automatic **assigning**, **ranking** and **classification**
- Problem of tossing sequence*
- Presence of recurring similar issues in historical maintenance calls

➤ **Ineffective root cause analysis and incident correlation**

- Need for deep fault localization and figure out dependencies among components and services

*Xie et al., Bug Triaging Based on Tossing Sequence Modeling. In **Journal of Computer Science and Technology 2019**

# Introduction and Motivation

❑ Capabilities of AI for Operating Systems (AIOps)†

**Prevention**

Forecast high-severity outages, future events, Alerting signals, Assessing system health

**Detection**

Detect abnormal conditions, Automated pattern discovery, Noise reduction in data

**Location**

Root cause analysis, Recurrent issues identification, unified topology and contextualization

**AIOps Abilities**

**Perception**

Data collection and ingestion, Data storage, Real-time monitoring, Querying data

**Interaction**

Human-computer intelligent interaction, Interactive analysis and collaboration

**Action**

Reactive Triage and routing, Prioritization of incidents, Set of Remediation actions

†**Remil et al.** AIOps Solutions for Incident Management: Technical Guidelines and A Comprehensive Literature Review, In **TOSEM 2023** [Under Submission]
***Dang et al.** AIOps: real-world challenges and research innovations. In **ICSE** 2019

A Data Mining Perspective on Explainable AIOps with Applications to Software Maintenance

# Introduction and Motivation

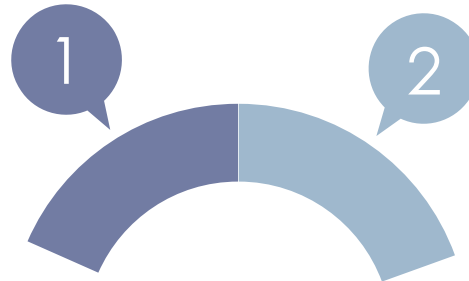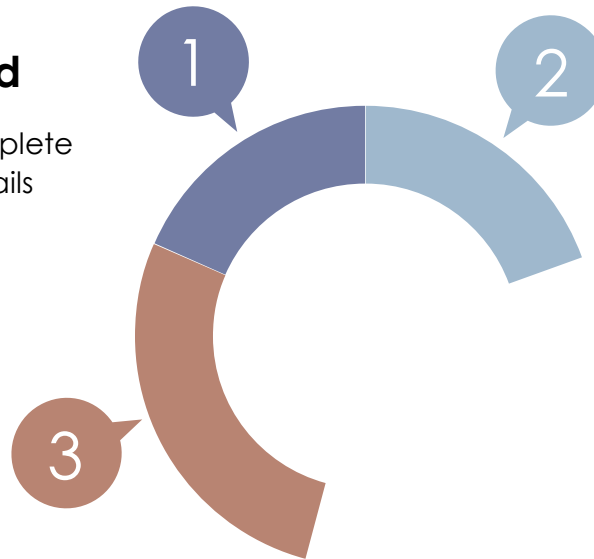❑ Research challenges of AIOps addressed in this thesis

# Introduction and Motivation

❑ Research challenges of AIOps addressed in this thesis

**Novel and Unstructured Field**

AIOps lacks unified terminology, complete taxonomy, desiderata, technical details
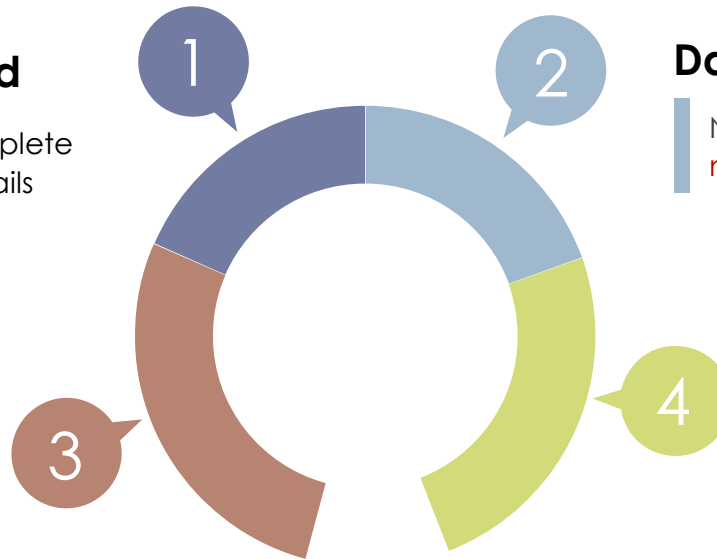
1

# Introduction and Motivation

❑ Research challenges of AIOps addressed in this thesis

**Novel and Unstructured Field**

AIOps lacks unified terminology, complete taxonomy, desiderata, technical details

**1**

**2**

**Data Requirements**

Noisy, unstructured, missing, unlabeled, non-homogeneous and complex data

# Introduction and Motivation

❑ Research challenges of AIOps addressed in this thesis

**Novel and Unstructured Field**

AIOps lacks unified terminology, complete taxonomy, desiderata, technical details

**Data Requirements**

Noisy, unstructured, missing, unlabeled, non-homogeneous and complex data

**Model Design**

Impractical supervised methods, overlooking descriptive models
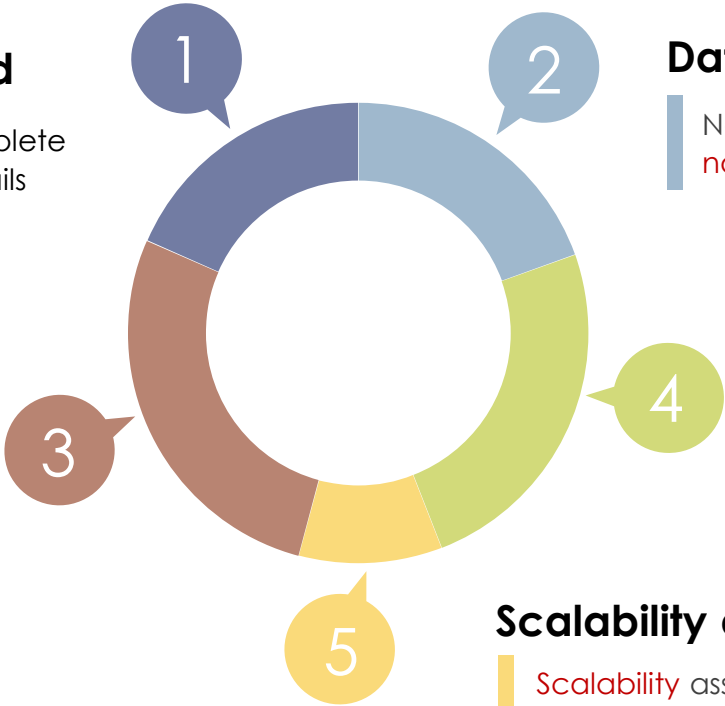
# Introduction and Motivation

❑ Research challenges of AIOps addressed in this thesis

**Novel and Unstructured Field**

AIOps lacks unified terminology, complete taxonomy, desiderata, technical details

**Data Requirements**

Noisy, unstructured, missing, unlabeled, non-homogeneous and complex data

**Model Design**

Impractical supervised methods, overlooking descriptive models

**Interpretability**

Best models are black box, transparency is preferred over performance

# Introduction and Motivation

❑ Research challenges of AIOps addressed in this thesis

**Novel and Unstructured Field**

AIOps lacks unified terminology, complete taxonomy, desiderata, technical details

**Data Requirements**

Noisy, unstructured, missing, unlabeled, non-homogeneous and complex data

**Model Design**

Impractical supervised methods, overlooking descriptive models

**Interpretability**

Best models are black box, transparency is preferred over performance

**Scalability and Robustness**

Scalability assessment often overlooked, temporal and in-context evaluation

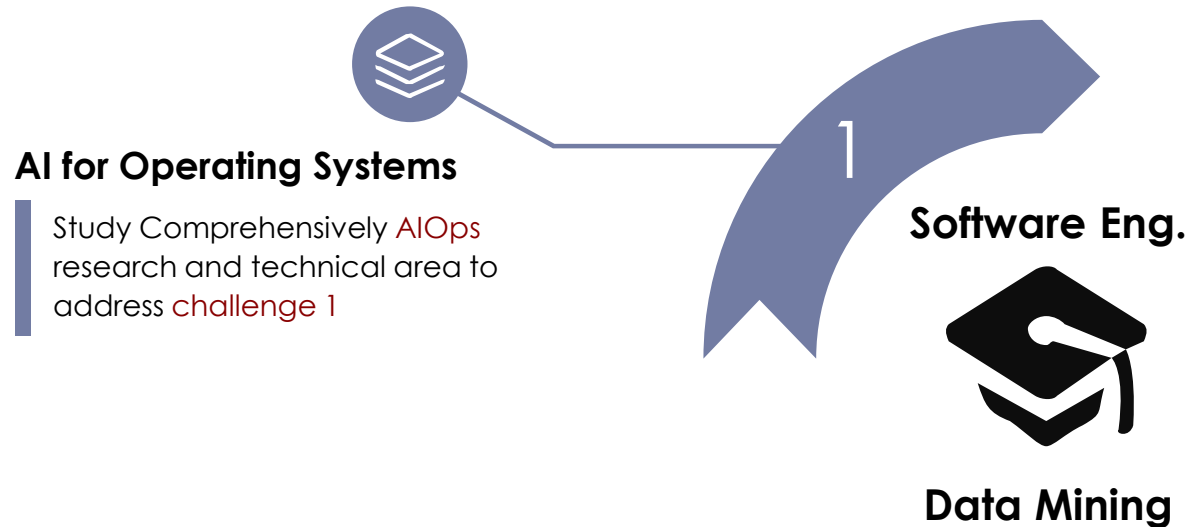# Introduction and Motivation

❑ Contributions and Key Research Areas

**Software Eng.**



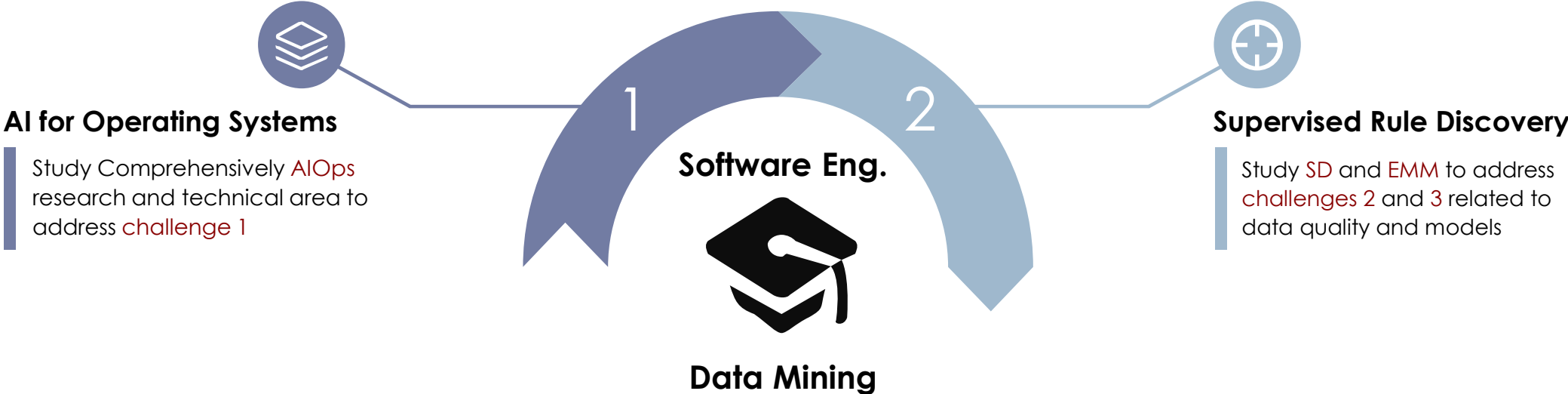**Data Mining**

# Introduction and Motivation

❑ Contributions and Key Research Areas

**AI for Operating Systems**

Study Comprehensively AIOps research and technical area to address challenge 1

1

**Software Eng.**

**Data Mining**

**Remil** et al. AIOps Solutions for Incident Management: Technical Guidelines and A Comprehensive Literature Review, In **TOSEM 2023 [Under revision, Core 2021, A\*]**
Bendimerad, **Remil** et al. On-premise Infrastructure for AIOps in a Software Editor SME: An Experience Report, In **ESEC/FSE 2023 [Published, Core 2021, A\*]**
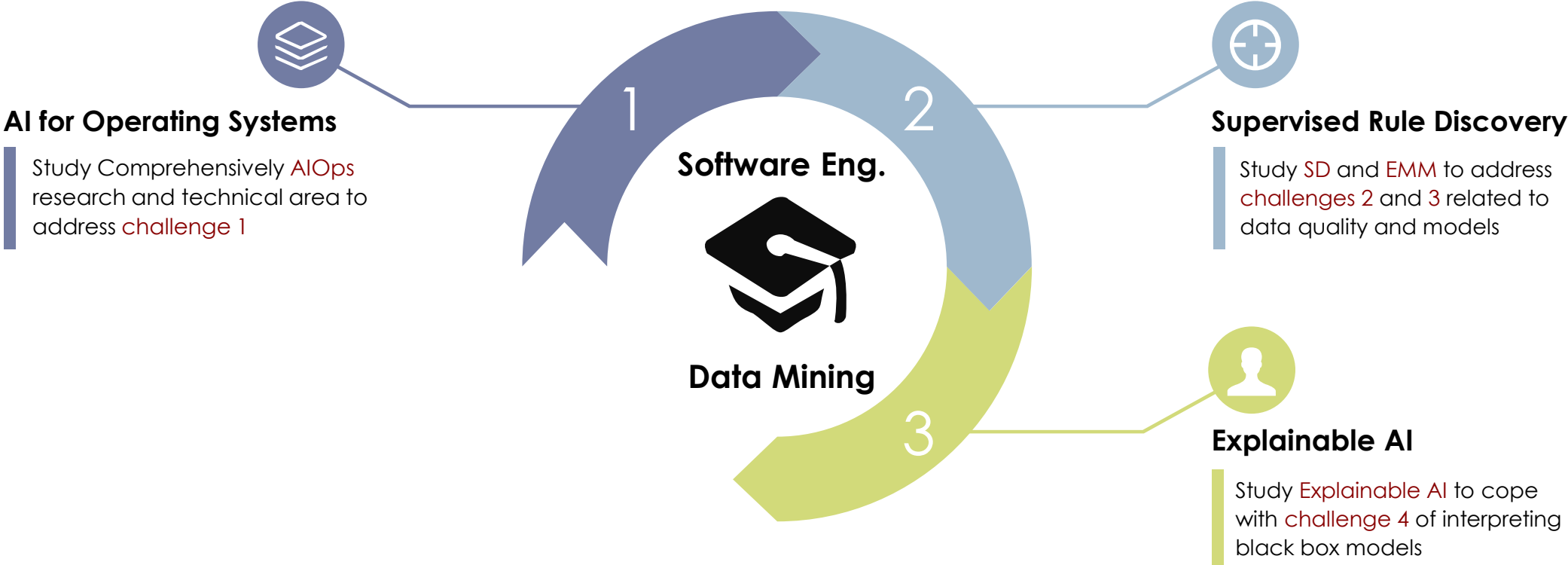
# Introduction and Motivation

❑ Contributions and Key Research Areas

**AI for Operating Systems**

Study Comprehensively AIOps research and technical area to address challenge 1

**1**

**Software Eng.**

**Data Mining**

**2**

**Supervised Rule Discovery**

Study SD and EMM to address challenges 2 and 3 related to data quality and models

**Remil** et al. What makes my queries slow: Subgroup Discovery for SQL Workload Analysis, In **ASE 2021 [Published**, **Core 2021, A\***]
**Remil** et al. Interpretable Summaries of Black Box Incident Triaging with Subgroup Discovery, In **DSAA 2021 [Published**, **Core 2021, A]**
**Remil** et al. Mining Java Memory Errors using Subjective Interesting Subgroups with Hierarchical Targets, In **ICDMW 2023 [Published, Workshop]**
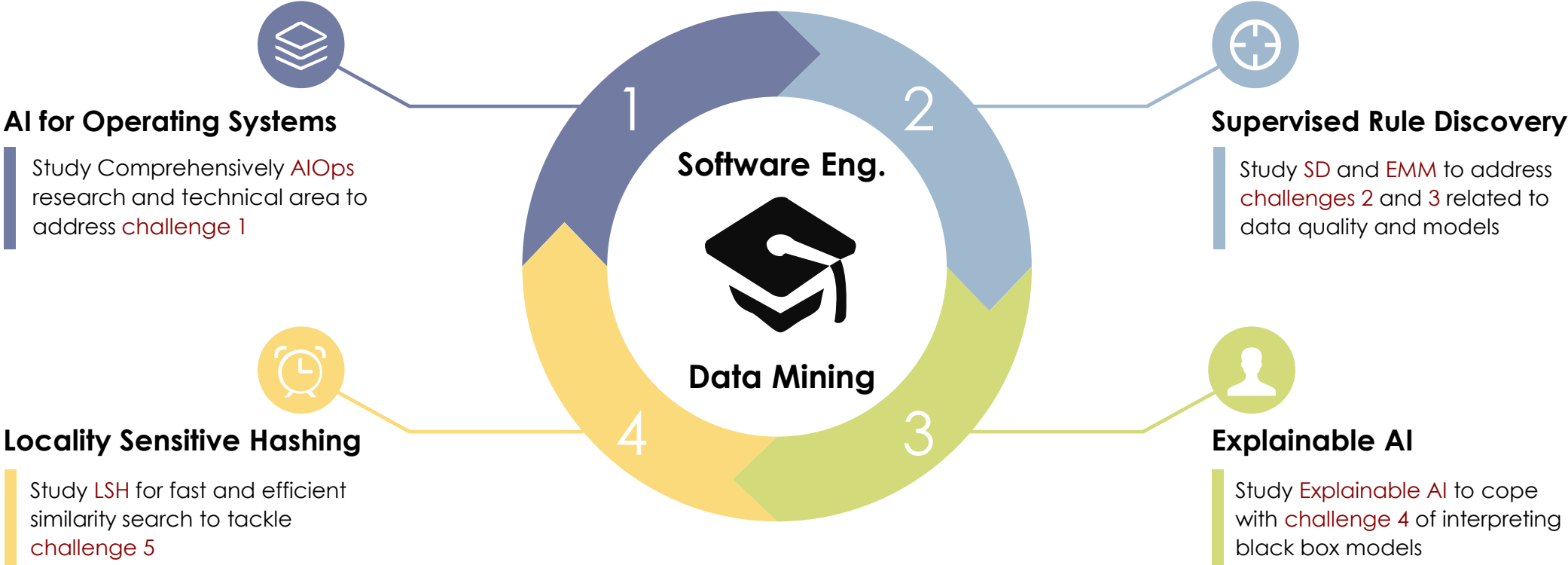
# Introduction and Motivation

❑ Contributions and Key Research Areas

**AI for Operating Systems**

Study Comprehensively AIOps research and technical area to address challenge 1

1  Software Eng.

2  **Supervised Rule Discovery**

Study SD and EMM to address challenges 2 and 3 related to data quality and models

Data Mining

3  **Explainable AI**

Study Explainable AI to cope with challenge 4 of interpreting black box models

**Remil** et al. Interpretable Summaries of Black Box Incident Triaging with Subgroup Discovery, In **DSAA 2021 [Published, Core 2021, A]**
**Remil** et al. Découverte de Sous-groupes Interprétables pour le Triage d'incidents, In **EGC 2022 [Published, National Conf]**
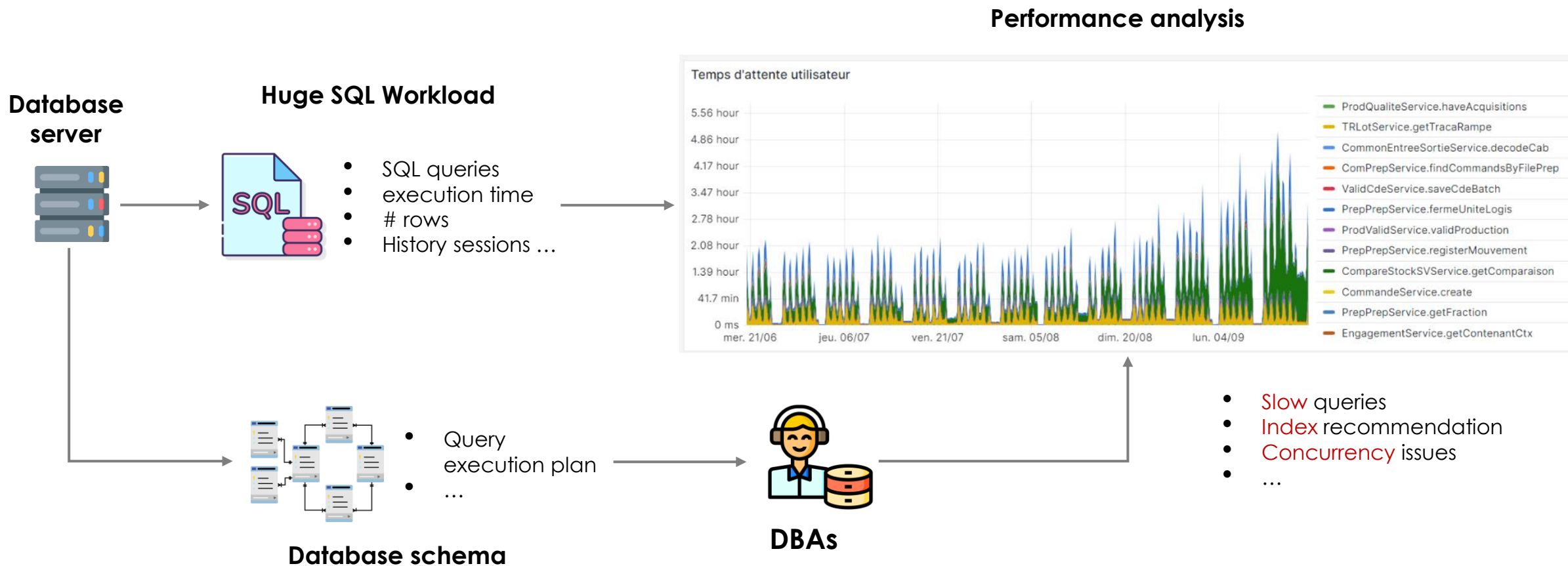
# Introduction and Motivation

❑ Contributions and Key Research Areas

**AI for Operating Systems**

Study Comprehensively AIOps research and technical area to address **challenge 1**

**Locality Sensitive Hashing**

Study LSH for fast and efficient similarity search to tackle **challenge 5**

**Software Eng.**

**Data Mining**

1

2

3

4

**Supervised Rule Discovery**

Study SD and EMM to address **challenges 2** and **3** related to data quality and models

**Explainable AI**

Study Explainable AI to cope with **challenge 4** of interpreting black box models

**Remil** et al. DeepLSH: Deep Locality-Sensitive Hash Learning for Fast and Efficient Near-Duplicate Crash Report Detection, In **ICSE 2024 [Published, Core 2021 A*]**

A Data Mining Perspective on Explainable AIOps with Applications to Software Maintenance

# Subgroup Discovery for SQL Workloads

**Performance analysis**

**Database server**

**Huge SQL Workload**

- SQL queries
- execution time
- # rows
- History sessions …

Temps d'attente utilisateur

5.56 hour
4.86 hour
4.17 hour
3.47 hour
2.78 hour
2.08 hour
1.39 hour
41.7 min
0 ms

mer. 21/06   jeu. 06/07   ven. 21/07   sam. 05/08   dim. 20/08   lun. 04/09

- ProdQualiteService.haveAcquisitions
- TRLotService.getTracaRampe
- CommonEntreeSortieService.decodeCab
- ComPrepService.findCommandsByFilePrep
- ValidCdeService.saveCdeBatch
- PrepPrepService.fermeUniteLogis
- ProdValidService.validProduction
- PrepPrepService.registerMouvement
- CompareStockSVService.getComparaison
- CommandeService.create
- PrepPrepService.getFraction
- EngagementService.getContenantCtx

- Query execution plan
- …

**Database schema**

**DBAs**

- Slow queries
- Index recommendation
- Concurrency issues
- …

# Subgroup Discovery for SQL Workloads

Need for a **generic** framework to analyse **batches** of SQL queries and bring answers to the question:
**How to characterize SQL queries that foster some properties of interest?**

# Subgroup Discovery for SQL Workloads

Need for a **generic** framework to analyse **batches** of SQL queries and bring answers to the question:
**How to characterize SQL queries that foster some properties of interest?**

**Illustrative example of SQL queries**

Pattern *P*: Predicate = verrou.date ∧ Db. Version = V2 ⟶ slow queries

| Predicates | | Topology | | … | Targets | |
|---|---|---|---|---|---|---|
| ik | date | Db. version | … | … | time | slow |
| 1 | **1** | **V2** | … | … | **14** | **1** |
| 0 | 1 | V1 | … | … | 2 | 0 |
| … | … | … | … | … | … | … |
| 1 | 0 | V2 | … | … | 3 | 0 |
| 0 | **1** | **V2** | | | **25** | **1** |



Overall data
**Pattern *P***

**SQL execution time probability**

**Atzmueller.** Subgroup Discovery, In **DAMI 2015**
**Wrobel.** An algorithm for multi-relational discovery of subgroups. In **PKDD 1997**

A Data Mining Perspective on Explainable AIOps with Applications to Software Maintenance

# Subgroup Discovery for SQL Workloads

**Subgroup Discovery building blocks**



symbol sets, numerical intervals, subgraphs, subsequences, … etc

**Langage Pattern $\mathcal{L}$**

**Enumeration Algorithm**

exhaustive, heuristic, by sampling, … etc

Selectors

Enumerate subgroups

**Dataset $\mathcal{D}$**

Evaluate subgroups

**Top $k$ subgroups**

Property of interest

numerical, binary, multiple complex attributes, … etc

**Target $T$**

f(x)

objective, subjective, semantic-based, constraints, … etc

**Interestingness $Q$**

# Subgroup Discovery for SQL Workloads

**SD building blocks for SQL Workload Analysis**

SQL queries

Parsing

Relevant data

Dataset $\mathcal{D}$

| Example of an SQL query |
|---|
| ```
SELECT m.ik
FROM model AS m
JOIN prod AS p
WHERE m.ik = p.ik
   AND m.uex = p1
   AND (m.uex in collection0
        OR m.ik in collection1)
   AND (m.dossier = p3
GROUP BY m.ik
HAVING (COUNT(DISTINCT p.ik) = p2)
AND (SUM(m.nbembal) = MAX (p.nbembal))
``` |

**Our parser***

| | |
|---|---|
| SELECT_model_ik | 1 |
| FROM_model<br>JOIN_prod | 1<br>1 |
| WHERE_model_ik<br>WHERE_model.uex<br>WHERE_model.dossier<br>WHERE_prod.ik | 2<br>1<br>1<br>1 |
| GROUPBY_ik | 1 |
| HAVING_prod_ik<br>HAVING_model.nbembal<br>HAVING_prod.nbembal | 1<br>1<br>1 |
| COUNT_prod.ik<br>SUM_model.nbembal<br>MAX_prod.nbembal | 1<br>1<br>1 |

*https://github.com/klahnakoski/mo-sql-parsing/pull/26

# Subgroup Discovery for SQL Workloads

**SD building blocks for SQL Workload Analysis**

SQL queries

Parsing

Relevant data

symbol sets,
numerical intervals,

**Langage Pattern $\mathcal{L}$**

Selectors

**Dataset $\mathcal{D}$**

Property of interest

numerical and
binary, attributes

**Target $T$**

Search Space $\mathcal{L}$ → Made of → **Selectors** → Conjunctive combinations → **Patterns**

- $sg(P) = ext(P) = \{c \in \mathcal{O} \, | P(o) = True\}$
- $P_{gen} \subset P_{spec} \Rightarrow sg(P_{gen}) \supseteq sg(P_{spec})$

DB version

V2

V3

Blocked
sessions

5   10   15   20   25

$\emptyset$

$sel_1 \quad sel_2 \quad \ldots \quad sel_d$

$sel_1 \wedge sel_2 \quad \ldots \quad sel_{d-1} \wedge sel_d$

$sel_1 \wedge sel_2 \wedge sel_3 \quad \ldots \quad sel_{d-2} \wedge sel_{d-1} \wedge sel_d$

$\bigwedge sel_i$

$P_{gen} : blockedSessions \in [15, 25]$
$P_{spec} : blockedSessions \in [15, 25] \wedge dbVersion = V3$

# Subgroup Discovery for SQL Workloads

**SD building blocks for SQL Workload Analysis**

SQL queries

Parsing

symbol sets,
numerical intervals,

Selectors

**Langage
Pattern** $\mathcal{L}$

Relevant data

**Dataset** $\mathcal{D}$

Property of interest

numerical attributes
e.g., runtime

Evaluate
subgroups

**Target** $T$

f(X)

**Interestingness** $Q$

Objective measures:
- **Exceptionality**
- **Generality**

Numerical measures used to evaluate the subgroup patterns*

1. **Mean**-based Measure (sensible to outliers)

$$q_{mean}^{\alpha} = i_P{}^{\alpha} \cdot (\mu_P - \mu_{\emptyset})$$

2. **Median**-based Measure (sensible to outliers)

$$q_{med}^{\alpha} = i_P{}^{\alpha} \cdot |Med_P - Med_{\emptyset}|$$

3. **T-Score** Measure (optimize the dispertion)

$$Tscore = i_P^{\frac{1}{2}} \cdot \frac{(\mu_P - \mu_{\emptyset})}{\sigma_P}$$

***Lemmerich.** Novel Techniques for Efficient and Effective Subgroup Discovery. **PhD thesis**, 2014

A Data Mining Perspective on Explainable AIOps with Applications to Software Maintenance

# Subgroup Discovery for SQL Workloads

**SD building blocks for SQL Workload Analysis**



SQL queries

Parsing

Relevant data

Dataset $\mathcal{D}$

symbol sets, numerical intervals,

Selectors

Langage Pattern $\mathcal{L}$

Enumerate subgroups

**Depth first search** Beam search

Evaluate subgroups

Refine

Property of interest

numerical and binary, attributes

Target $\mathcal{T}$

f(X)

Interestingness $Q$

Objective measures: Exceptionality Generality

**Enumeration Algorithm**

Empty pattern

Search Space

Max support

Pruning

**Lemmerich et al.**, Fast exhaustive subgroup discovery with numerical target concepts. In **DAMI 2016**

A Data Mining Perspective on Explainable AIOps with Applications to Software Maintenance

# Subgroup Discovery for SQL Workloads

**Results on a large workload of Hibernate queries made available by Infologic**

| ID | Target | Measure | Subgroup patterns | Size | Quality |
|---|---|---|---|---|---|
| D1 | time | Median | $(P_1)$ : WHERE_stocks.gestion.modele.lot.prod.ref.auditinfo.etat $\geqslant 1$ | 8 | $161 \cdot \text{q\_med}(P_\varnothing)$ |
| | | | $(P_2)$ : FROM_ventes.cumuls.modele.cumulmultiple $\geqslant 1$ | 451 | $21 \cdot \text{q\_med}(P_\varnothing)$ |
| | | | $(P_3)$ : WHERE_ventes.cumuls.modele.cumulmultiple.valzvcliX $\geqslant 1$ | 45 | $21 \cdot \text{q\_med}(P_\varnothing)$ |
| | | | $(P_4)$ : WHERE_.ventes.cumuls.modele.cumulmultiple.valzvartX $\geqslant 1$ | 45 | $21 \cdot \text{q\_med}(P_\varnothing)$ |
| D2 | slow $\tau_{P_\varnothing} \simeq 0.6$ | Lift | $(P_5)$ : GROUPBY_stocks.gestion.modele.mvtrealise.refexterne $\geqslant 1$ | 131 | $\tau_P = 1$ |
| | | | $(P_6)$ : serverName $=$ ServerX $\wedge$ systemI/O $> 50$ | 38 | $\tau_P = 1$ |
| | | WRAcc | $(P_7)$ : WHERE_stocks.gestion.modele.mvtrealise.etatsynchro $\geqslant 1 \wedge$ jdbcMax $< 200$ | 20668 | $\tau_P \simeq 0.99$ |
| | | | $(P_8)$ : WHERE_stocks.gestion.modele.mvtrealise.auditinfo.datcre $\geqslant 1 \wedge$ dbVersion $= 2.3$ | 20675 | $\tau_P \simeq 0.99$ |
| | | | $(P_9)$ : manyActiveSessions $=$ Alarm | 44 | $\tau_P \simeq 93\%$ |

# Enhancing Duplicate Crash Report Retrieval

**Problem of deduplication\***

A bug reported for the **ACTI** service

But the bug is generic and related to a **web** feature



à 11:50

: à marquer comme traiter puisque tu as fait une correction spécifique ici, on analysera le cas général dans WEBBUG-6883.

**Near-duplicates**

≠

**\*Jiang et al.** Igor: Crash Deduplication Through Root-Cause Clustering, In **CCS 2021**

A Data Mining Perspective on Explainable AIOps with Applications to Software Maintenance

# Enhancing Duplicate Crash Report Retrieval

**Similarity measures** for stack trace comparison embedded in **Clustering** algorithms

**Complex** similarity measures based

**Computational** Complexity is very **costly**

Measures embedded in **clustering** with several **issues**

It should be handle as **Nearest Neighbours Search** problem



| Call Stack $C_1$ | Call Stack $C_2$ |
|---|---|

Distance to Crash Point for $f_6'$:6

Distance to Crash Point for $f_4$:4

Alignment Offset between $f_4$ and $f_6'$ : 2

Matched Frame
Other Frame

[**Dang et al.**, in **ICSE 2012**]*

***Dang et al.** ReBucket: A Method for Clustering Duplicate Crash Reports Based on Call Stack Similarity. In **ICSE 2012**
†**Wu et al.** CrashLocator: Locating Crashing Faults Based on Crash Stacks. In **ISSTA 2013**
‡**Moroo et al.** Reranking-based Crash Report Deduplication. **SEKE 2017**

**Contribution**

**Learn a family of hash functions with a constrained hashing Siamese neural network**

# Enhancing Duplicate Crash Report Retrieval

**Does the model manage to converge to the LSH property?**

# Enhancing Duplicate Crash Report Retrieval

**Is the model fast enough compared to linear scans?**

| Similarity Measure | Runtime (~ Seconds) | | | | |
|---|---|---|---|---|---|
| | k-NN | CNNH+LSH | DeepLSH | MinHash | SimHash |
| Jaccard | 258 | 30 | 26 | 57 | - |
| Cosine | 8288 | 15 | 14 | - | 3 |
| TF-IDF | 8510 | 16 | 15 | - | 4 |
| Edit Distance | 4911 | 29 | 29 | - | - |
| PDM | 10047 | 16 | 16 | - | - |
| Brodie | Limit | 27 | 27 | - | - |
| DURFEX | 12160 | 26 | 24 | - | - |
| Lerch | 3118 | 24 | 24 | - | - |
| Moroo | 15253 | 25 | 25 | - | - |
| TraceSim | 13050 | 30 | 30 | - | - |

# End

# Thanks for your attention

# End

# Q/A?

A Data Mining Perspective on Explainable AIOps with Applications to Software Maintenance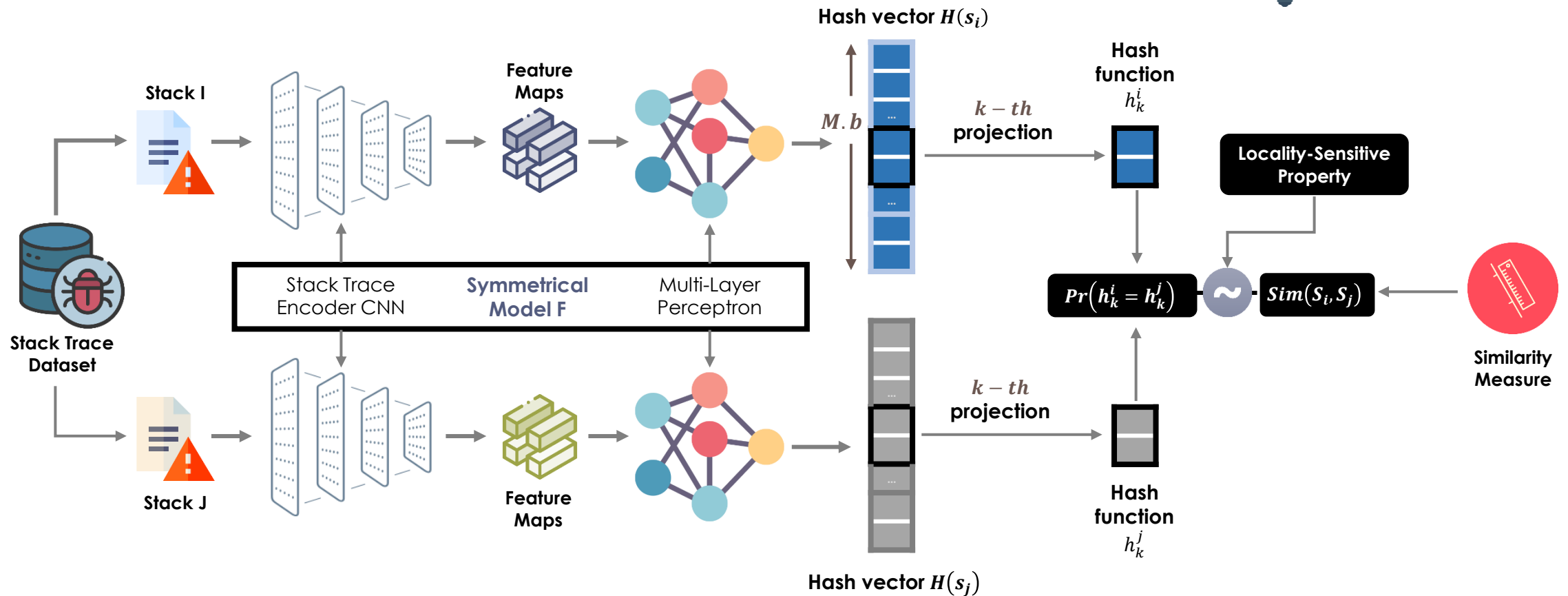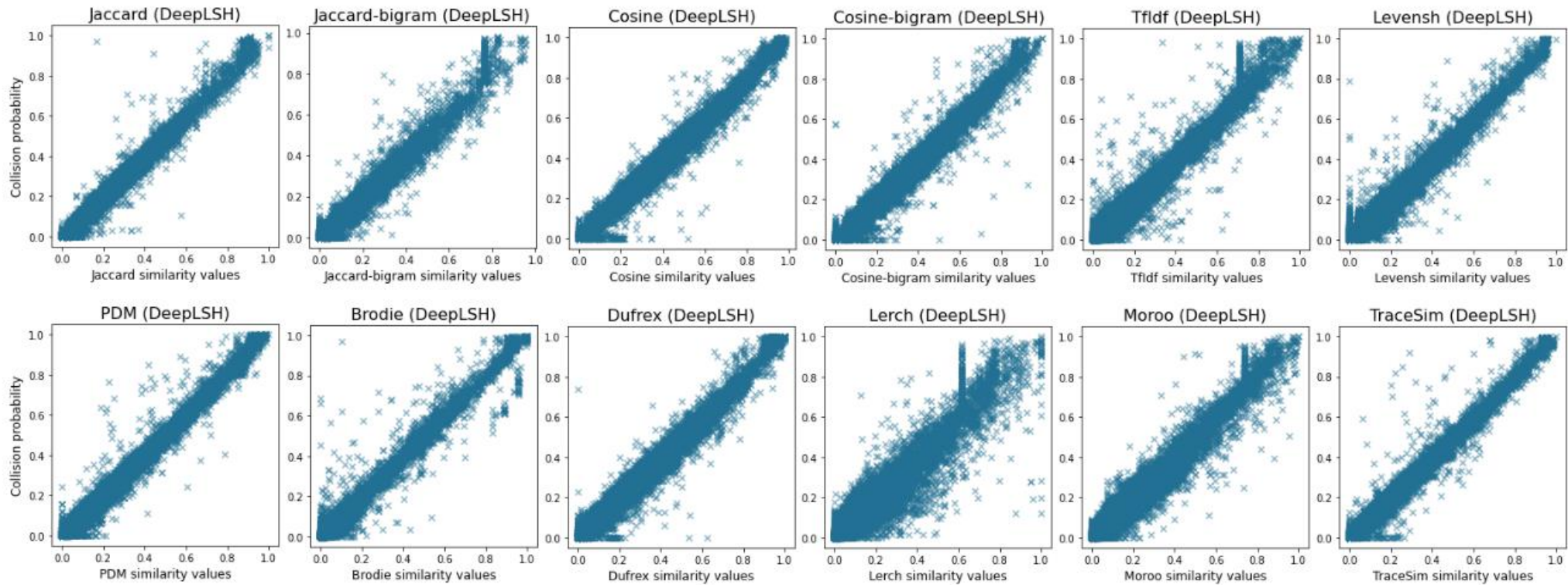