



The 12th Symposium on Conformal and Probabilistic Prediction with Applications
Limassol, Cyprus
September 12th – 15th 2023

Enterprise Disk Drive Scrubbing Based on Mondrian Conformal Predictors

Rahul Vishwakarma⁺

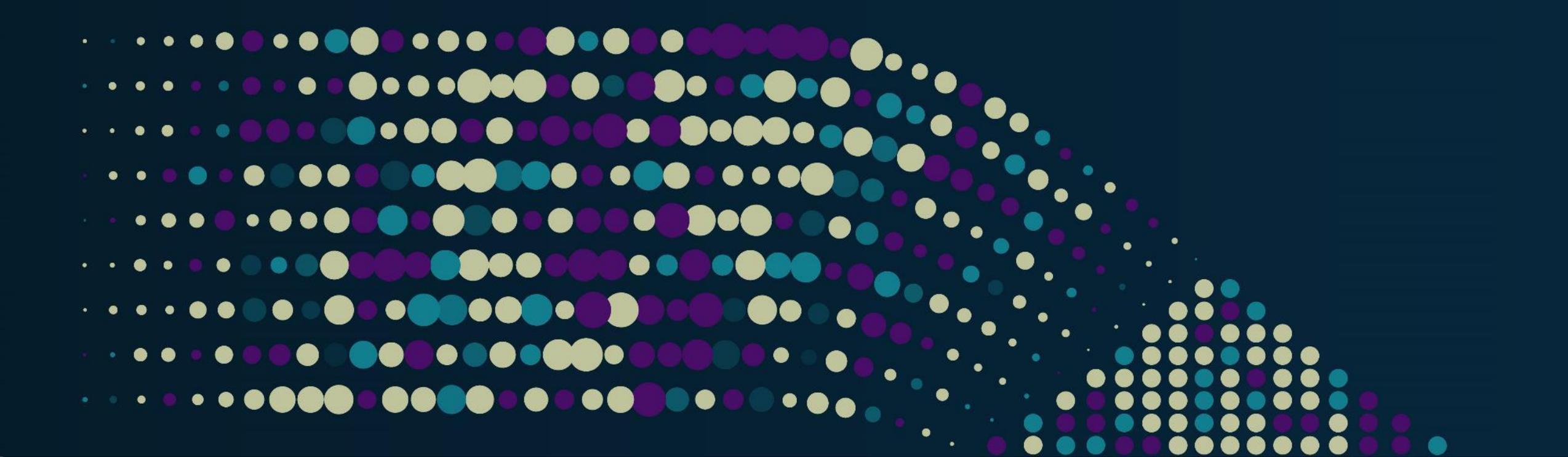
Jinha Hwang⁺

Soundouss Messoudi^o

Ava Hedayatipour⁺

California State University Long Beach, USA ⁺

HEUDIASYC - UMR CNRS 7253, Université de Technologie de Compiègne, France ^o



Introduction

Introduction

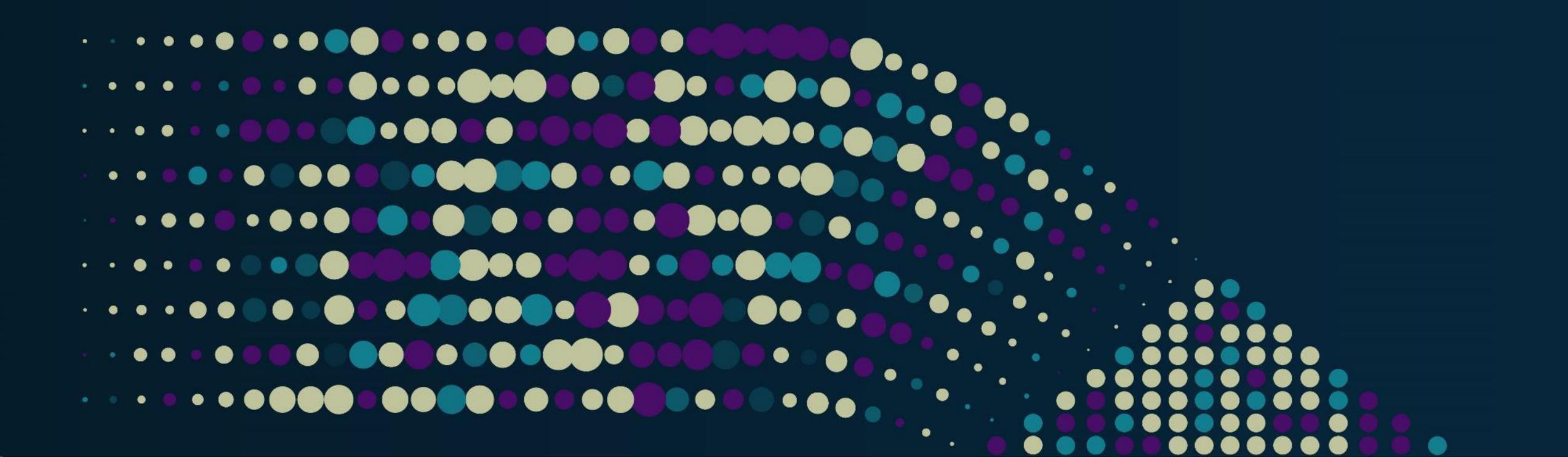
- Data centers are intricate ecosystems with diverse storage devices spread across multiple locations.
- Managing business becomes challenging due to the complexity of each individual storage component.
- Inadequate management of storage components can lead to data loss and business disruption.
- Proactive monitoring of storage component health is essential, with e.g. **disk drive scrubbing**.



Data center illustration



Which disks to scrub ? When to scrub ?



Disk Scrubbing

What is Disk Scrubbing ?

- Error correction technique :
 - Periodically checks for inconsistencies in data, then corrects them.
 - Applied to **disk arrays**, memory, file systems, ...

- Targeted problems :
 - Data integrity issues.
 - Degraded system performance.
 - Reduced drive lifespan.
 - Regulatory and compliance issues.



Hard Disk Drive



Solid State Drive

Disk Scrubbing challenges

■ Resource Consumption

- Data scrubbing is resource-intensive, consuming CPU, memory, and storage I/O.
- Frequent scrubbing strains system resources, impacting system performance.

■ Drive Lifespan Impact

- Frequent scrubbing can accelerate wear and tear on disk drives.
- Counterproductive if it accelerates drive degradation instead of preventing it.

■ Workload Impact

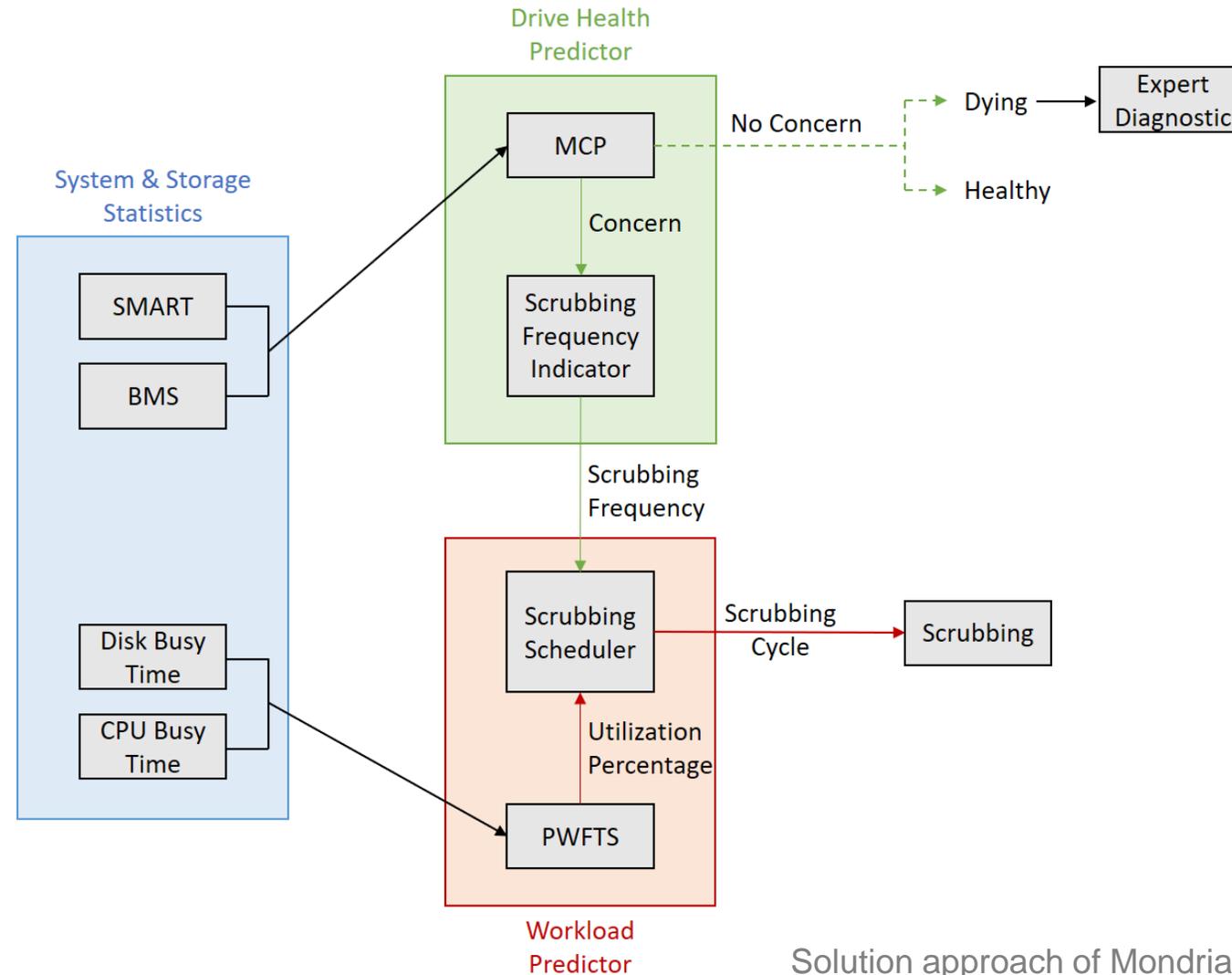
- Excessive scrubbing can affect the performance of concurrent workloads.
- May lead to slowdowns or interruptions in critical operations.

Quest for a fine balance between **scrubbing frequency** and **system reliability**.



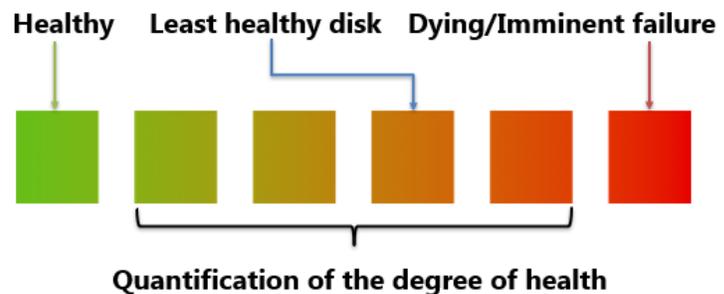
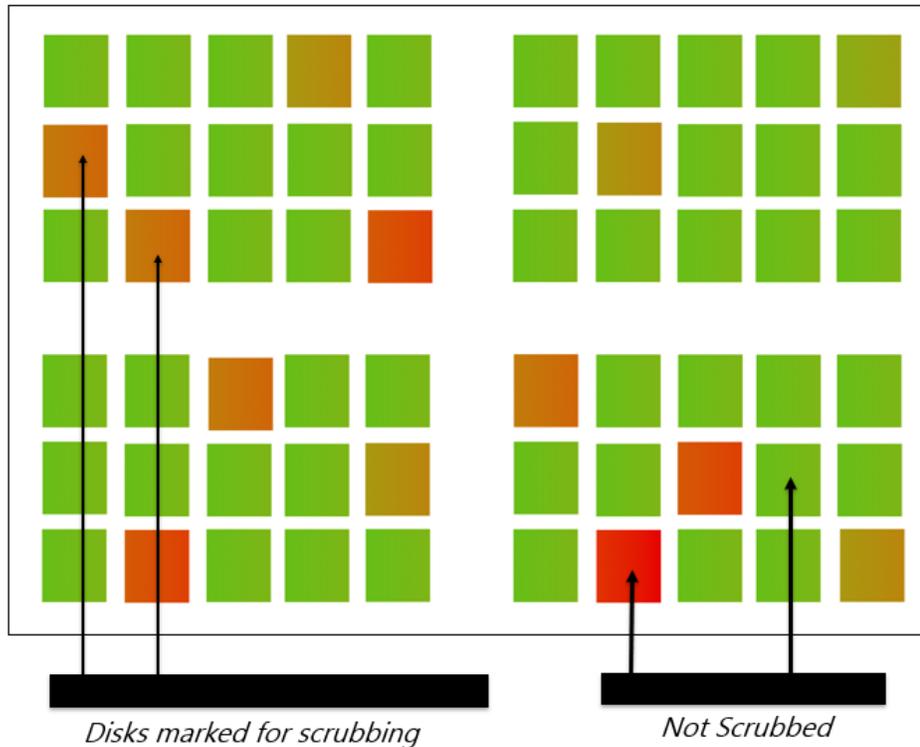
Approach : Mondrian Conformal Disk Drive Scrubbing

Solution Architecture



Solution approach of Mondrian conformal disk drive scrubbing

Which disks to scrub : Drive Health Predictor



- **Selective scrubbing**

- Eliminate “no concern” disks (healthy / dying).
- Identify the concerned disks.

- **Degree of health**

- Good, medium or poor.

- **Scrubbing cycle**

- Low, medium, high (respectively).



Mondrian Conformal Prediction

Mondrian Conformal Prediction

■ Why ?

- In most cases, drive health prediction is highly imbalanced.

$$p_{n+1}^{C_k} = \frac{|\{i \in 1, \dots, q : y_i = C_k, \alpha_{n+1}^{C_k} \leq \alpha_i\}|}{|\{i \in 1, \dots, q : y_i = C_k\}|}$$

■ How ?

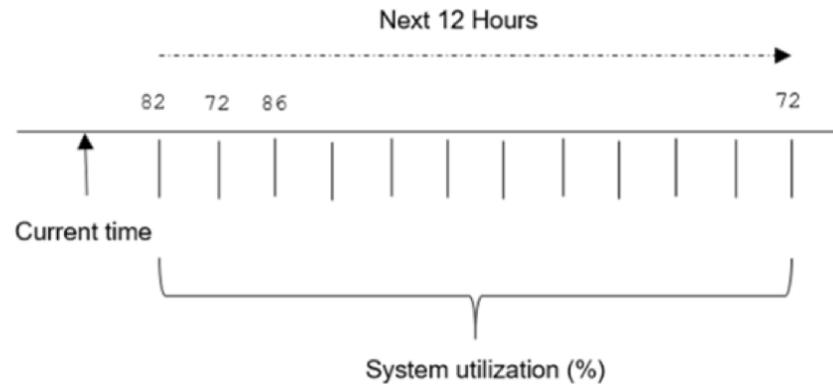
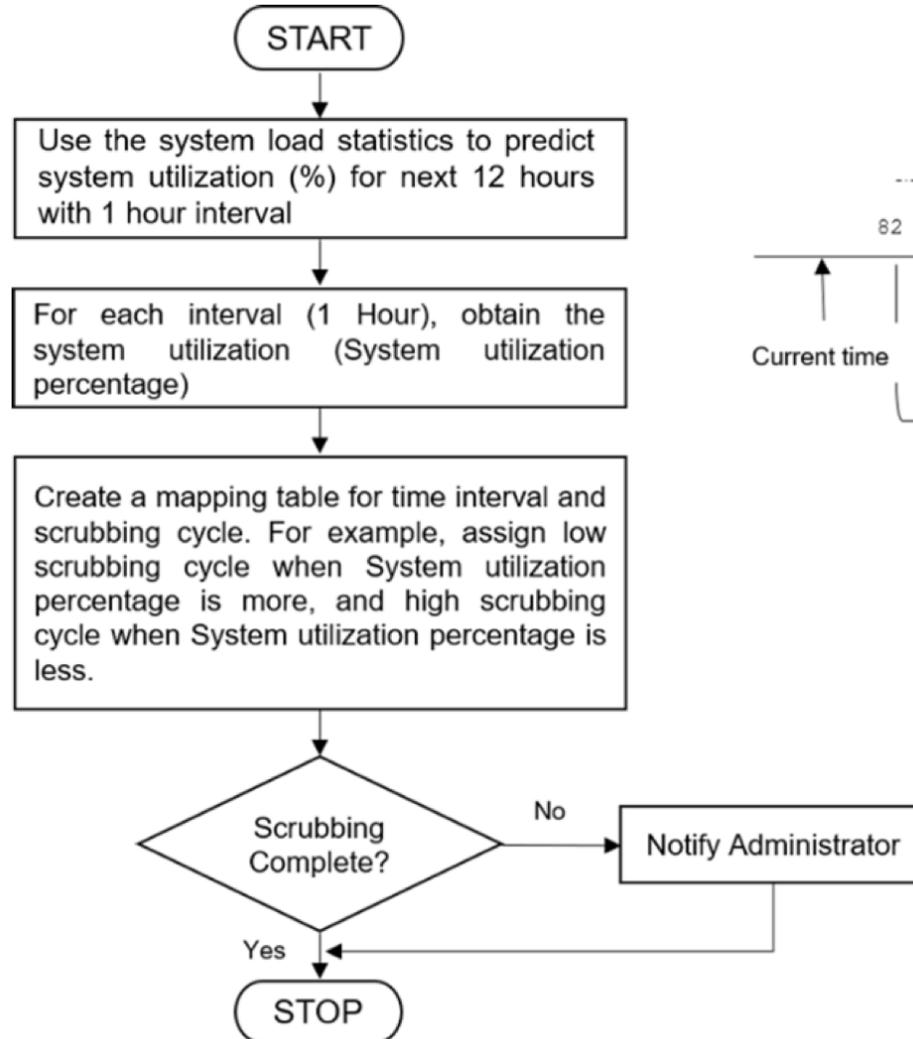
- Binary classification.
- Confidence score as a health score.

$$\text{Confidence}(x) = \sup\{1 - \epsilon : |\Gamma_\epsilon(x)| \leq 1\}.$$

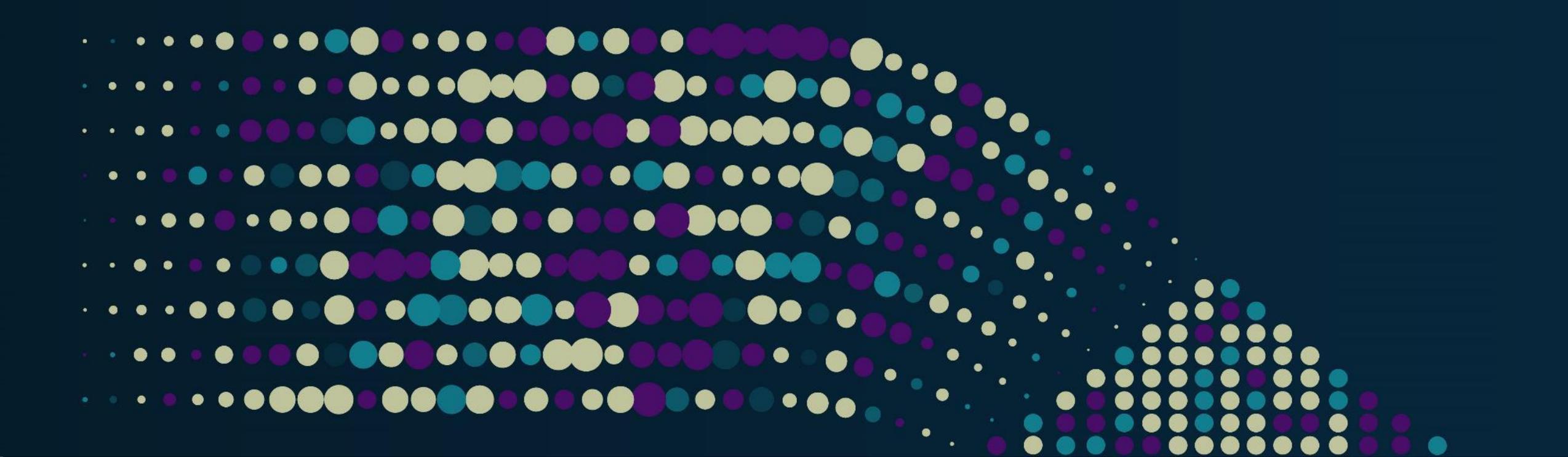
Disk health	Scrubbing frequency	Health score
Best	LOW	[95%, 99[
Medium	MEDIUM	[80%, 95[
Poor	HIGH	< 80%

Mapping of the disk health with the scrubbing frequency based on health score.

When to scrub : Workload Predictor



Flowchart of the workload predictor using the PWFTS algorithm



Experiments

Open-source Baidu dataset

- Seagate ST31000524NS* (HDD).

Column No	Feature	Description
1	Index of the disk serial number	Ranging from 1 to 23395
2	Value of SMART ID #1	Raw Read Error Rate
3	Value of SMART ID #3	Spin Up Time
4	Value of SMART ID #5	Reallocated Sectors Count
5	Value of SMART ID #7	Seek Error Rate
6	Value of SMART ID #9	Power On Hours
7	Value of SMART ID #187	Reported Uncorrectable Errors
8	Value of SMART ID #189	High Fly Writes
9	Value of SMART ID #194	Temperature Celsius
10	Value of SMART ID #195	Hardware ECC Recovered
11	Value of SMART ID #197	Current Pending Sector Count
12	Raw Value of SMART ID #5	Reallocated Sectors Count
13	Raw Value of SMART ID #197	Current Pending Sector Count
14	Class label of the disk	0 for failed and 1 for functional

Features' description for the Open-source Baidu dataset.

* Hdds dataset (baidu inc.), Jan 2023. <https://www.kaggle.com/datasets/drtycoon/hdds-dataset-baidu-inc>.

Results

- MCP identifies more disks of the minority class, but with a decrease in the number of disks correctly classified as healthy...

		kNN		MCP		
		Predicted				
			0	1	0	1
Actual	0		51314	547	51669	537
	1		975	296689	28703	268616

Comparison of confusion matrix results for disk drive classification using kNN and MCP.

- But by considering the health score:

Health score	0.998 - 0.9985	0.9985 - 0.999	0.999 - 0.9995	0.9995 and 1
Disk count	16468	62928	63087	126244

The number of relatively healthy drives based on the health score intervals.



Only 22.7% of the number of disks is scrubbed.

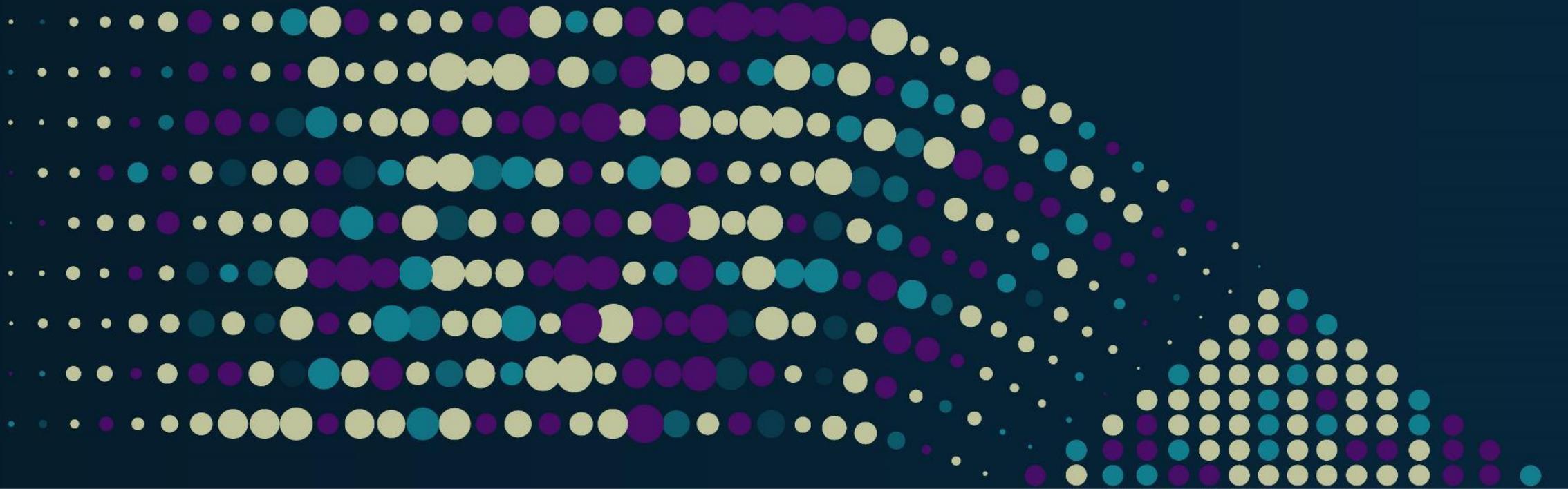
Real life datasets

- **Data 1:** 2015-01-01—2019-07-01, Hitachi drives
- **Data 2:** 2015-01-01—2019-07-01, 3T ST330006CLAR3000

Confidence	Data_set_one	Data_set_two
>=99%	14732	8385
>=98%, <99%	2401	144
>=97% , <98%	5688	55
>=96%, <97%	2257	20
>=95%, <96%	1058	19
>=94%, <95%	680	9
>=93%, <94%	422	4
>=92%, <93%	94	2
>=91%, <92%	0	0
>=90%, <91%	0	0
<90%	0	0
total	27332	8638

Experiments on real life data centers.

➔ Only 12600 (data 1) and 253 (data 2) disks are scrubbed, with different frequencies of scrubbing.



Conclusions and perspectives

Conclusions

■ **New Scrubbing Approach**

- Introduces a fine-grained method for selecting drives to be scrubbed.
- Enhances the existing failure analysis engine with an algorithm-agnostic Mondrian conformal predictor.

■ **Confidence-Based Health Ranking**

- Translates prediction confidence into a disk health ranking mechanism.
- Transforms this ranking into a scrubbing frequency depending on the user's preferences.

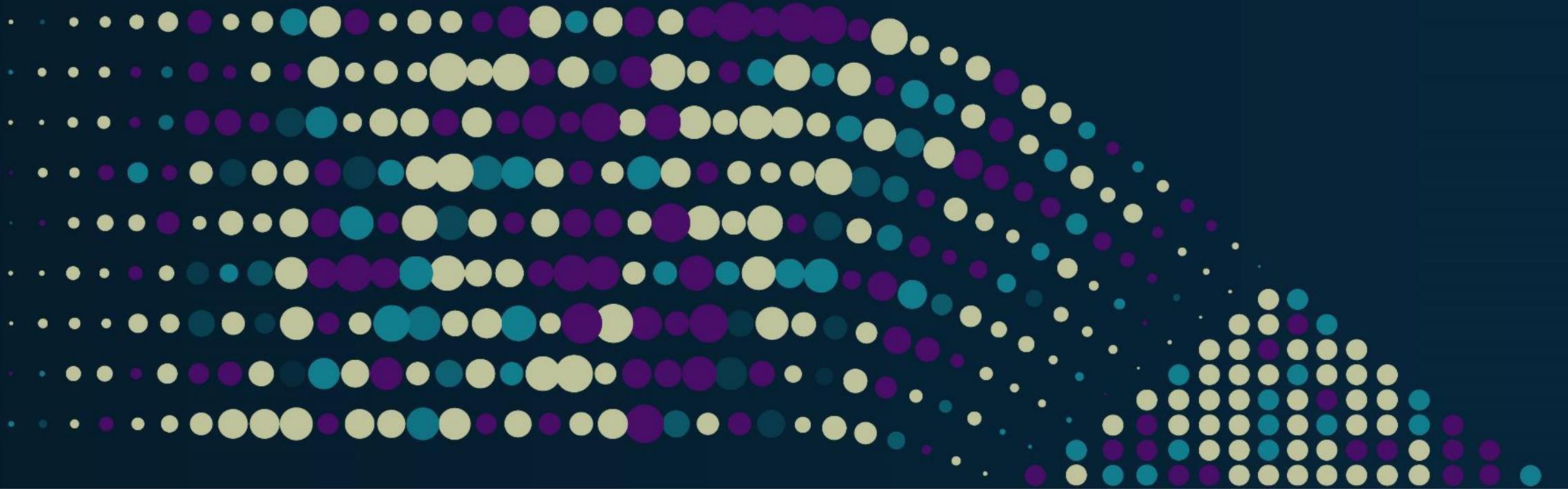
■ **Optimized Scrubbing Schedule**

- Uses n-step ahead system load prediction.
- Optimizes disk scrubbing for improved efficiency.

■ **Energy saving and reduced carbon footprint for data centers**

Perspectives

- Adapt the health drive predictor over time and system utilization.
- Improve the current approach by exploring other non-conformity measures.
- Apply the approach in large-scale data centers, preferably in real-life cases.
- Consider Venn-Abers predictors with calibrated probabilities for predictions.



Thank you for your attention