

Ines Arous<sup>1</sup>, Mourad Khayati<sup>1</sup>, Philippe Cudré-Mauroux<sup>1</sup>  
Ying Zhang<sup>2</sup>, Martin Kersten<sup>2</sup>, Svetlin Stalinlov<sup>2</sup>



<sup>1</sup>{firstname.lastname}@unifr.ch

<sup>2</sup>{firstname.lastname}@monetdbolutions.com

## GOAL AND CONTRIBUTIONS

**Motivation:** Real-world time series (sensor) data often contain missing values. Missing values are harmful to upper-level time series analytics.

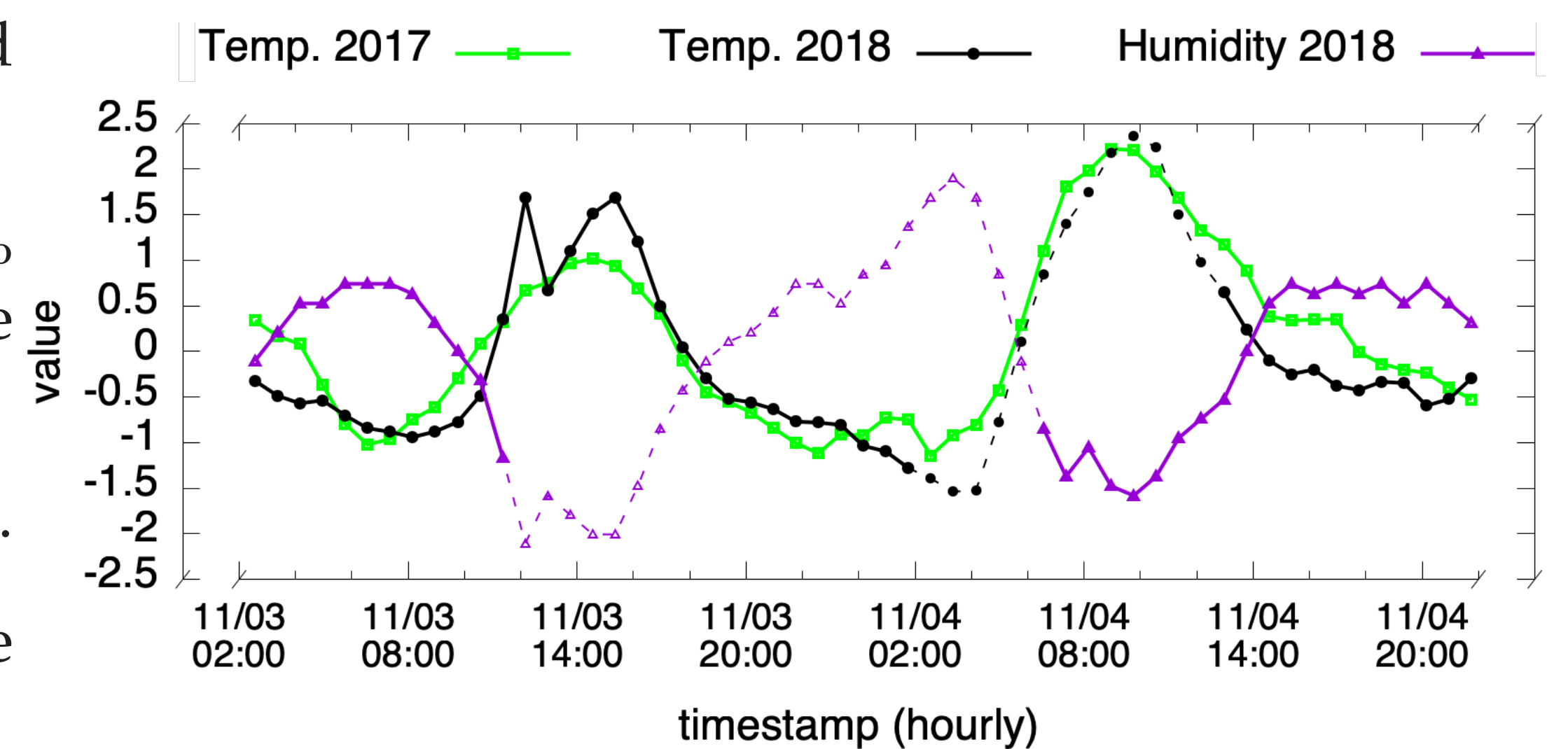
**Goal:** An efficient and accurate DB system to recover missing values in time series.

**Contribution:** A system called RECOVDB with:

- Parameter-free recovery
- Correlation-aware recovery
- Full-fledged DBMS (MonetDB) support

## RECOVERY CHALLENGES

- Long and (linearly) correlated time series.
- Large missing blocks (up to 80% of consecutive observations are missing per time series).
- Multiple incomplete time series.
- Integrate the recovery in the MonetDB system.



## RECOVERY GUI

### Menu

Data:  Raw  Z-Score  Min-Max

### Hydrology: Discharge of Swiss rivers

Source: Federal Office for the Environment FOEN  
Unit: m3/s

### Recovery of missing values

Recover missing values for: ▾

- Appenzell
- Halden
- Jonschwil
- Liestal
- Moutier
- Rheinhalle
- Wilier

Threshold epsilon for CD:

0.01

Drop values by:

40%

Apply

use PHP

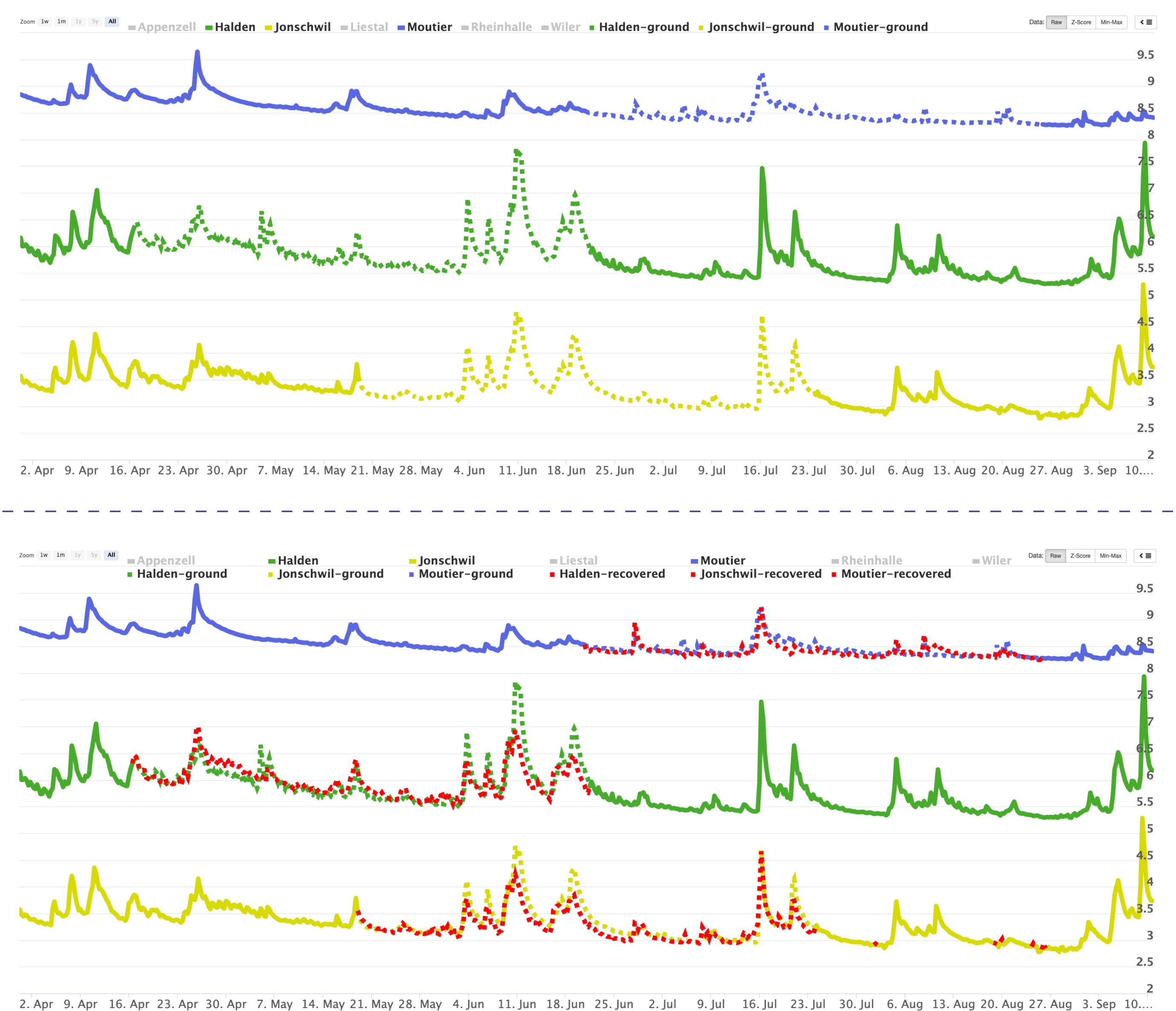
Recover

1. Data normalization options.
2. Select one or multiple TS from which to drop blocks of missing values.
3. Select a recovery termination threshold.
4. Select percentage of additional missing values.
5. Use PHP-based recovery instead of UDF.
6. Recover the missing blocks with the selected setup.

- The tool is accessible via [revival.exascale.info](http://revival.exascale.info)
- The tool can be easily extended with new datasets

### Recovery Example

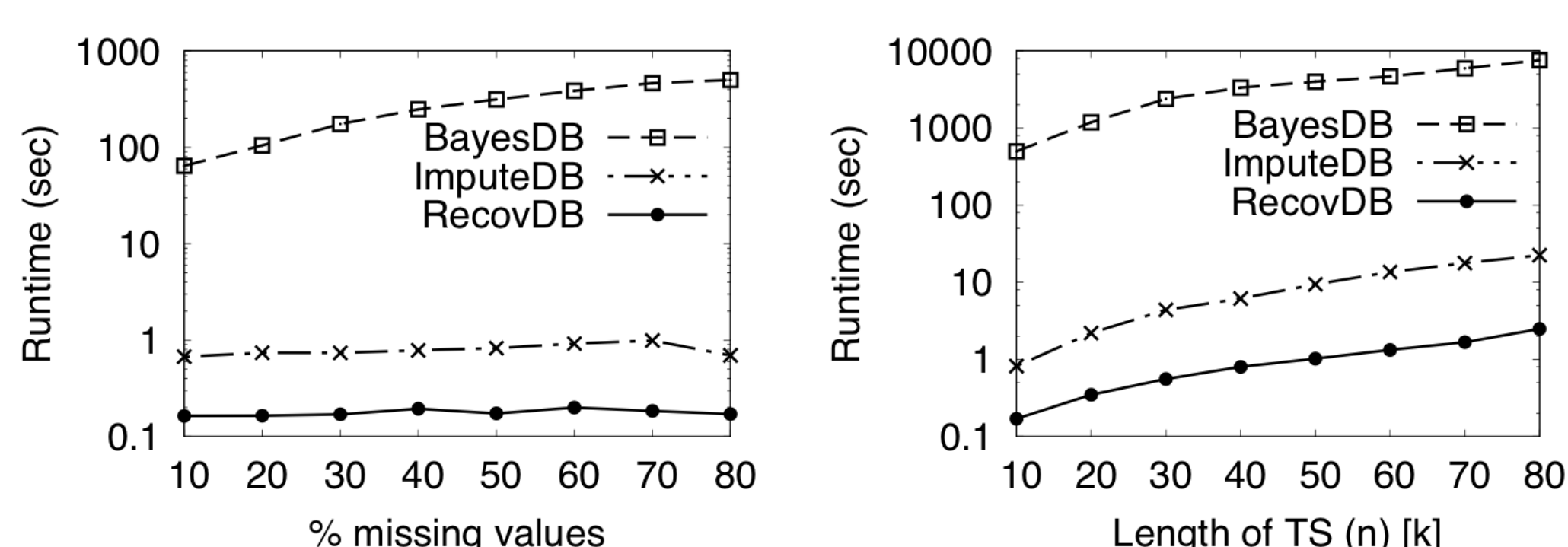
- Missing values are denoted by dashed lines and recovered values are denoted by red lines.
- RECOVDB leverages the inter-time series correlation to accurately recover multiple TS in one go.



## EMPIRICAL EVALUATION

### Efficiency

- RECOVDB is up to 10000x and 10x faster than BayesDB<sup>a</sup> and ImputeDB<sup>b</sup> respectively.
- To achieve this performance, RECOVDB exploits the analytical power of MonetDB to handle the data management and pre-/post-processing.

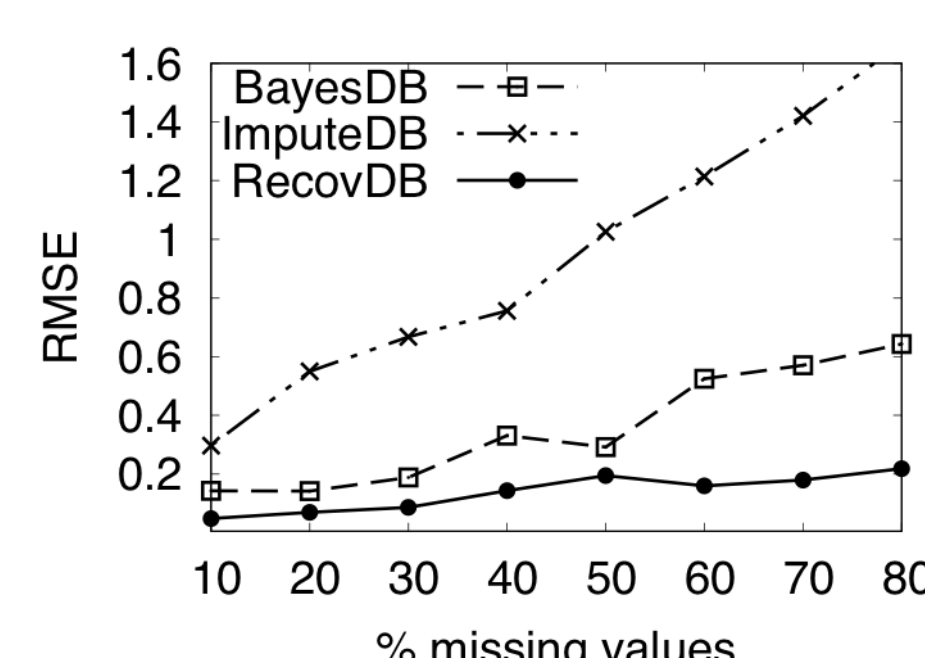


<sup>a</sup>Saad, F. and Mansinghka, V.K. *A probabilistic programming approach to probabilistic data analysis*. NIPS, 2016

<sup>b</sup>Cambronero, J., Feser, J.K., Smith, M.J. and Madden, S. *Query optimization for dynamic imputation* VLDB, 2017

### Accuracy

- RECOVDB outperforms the STOA using RMSE, MSE and MAE metrics.
- RECOVDB is up to 66% and 87% more accurate than BayesDB and ImputeDB, respectively.
- The accuracy of RECOVDB is steady with increasing % of missing values.



## DEMO SCENARIOS

**Scenario 1** Recover multiple incomplete time series at one.

**Scenario 2** Increase the size of the missing block and of the data.

**Scenario 3** Compare RECOVDB against STOA recovery DB systems (i.e., BayesDB and ImputeDB).

## CONCLUSIONS

- We present RECOVDB which recovers large missing blocks in multiple time series.
- RECOVDB leverages the correlation across time series during the recovery.
- RECOVDB outperforms STOA in both efficiency and accuracy when increasing i) the length/number of time series and ii) the size of the missing blocks.